

The language of civil engineering research articles ~ a corpus-based approach

Alex Gilmore
Department of Civil Engineering
University of Tokyo, Japan

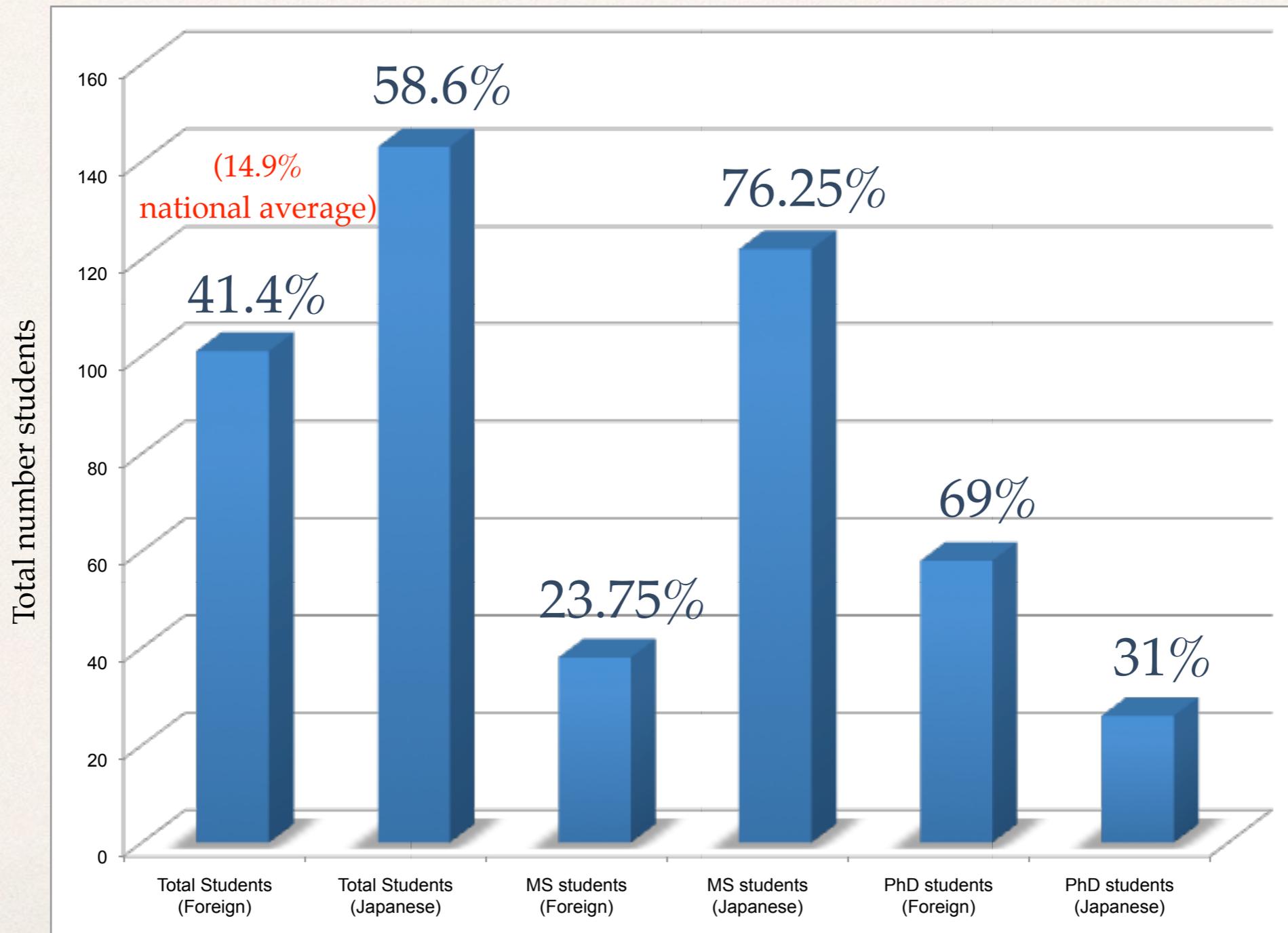


Department of Civil Engineering, University of Tokyo

Department of Civil Engineering

1. Infrastructure Development & Construction Management
2. Regional Planning & Surveying
3. Transportation Engineering & Planning
4. River & Environmental Engineering
5. Coastal & Ocean Engineering
6. Hydrology & Water Resources Engineering
7. Geotechnical Engineering
8. Concrete & Construction Engineering
9. Earthquake & Disaster Mitigation Engineering
10. Mechanics & Structures
11. International Projects

Postgraduate student population - Department of Civil Engineering, University of Tokyo



Incoming students: Where are they from?



Why create a specialised corpus?

- * Large variation between different academic disciplines in terms of word frequencies, collocational patterns & rhetorical moves: e.g. 4-word lexical bundles from fields of Biology, Electrical Engineering, Applied Linguistics & Business Studies >50% unique (Hyland 2008)
- * Specialised corpora a good starting point for design of ESP materials - an area where publishers are unwilling to invest resources due to the restricted audience

SCCERA (Specialized Corpus of Civil Engineering Research Articles)

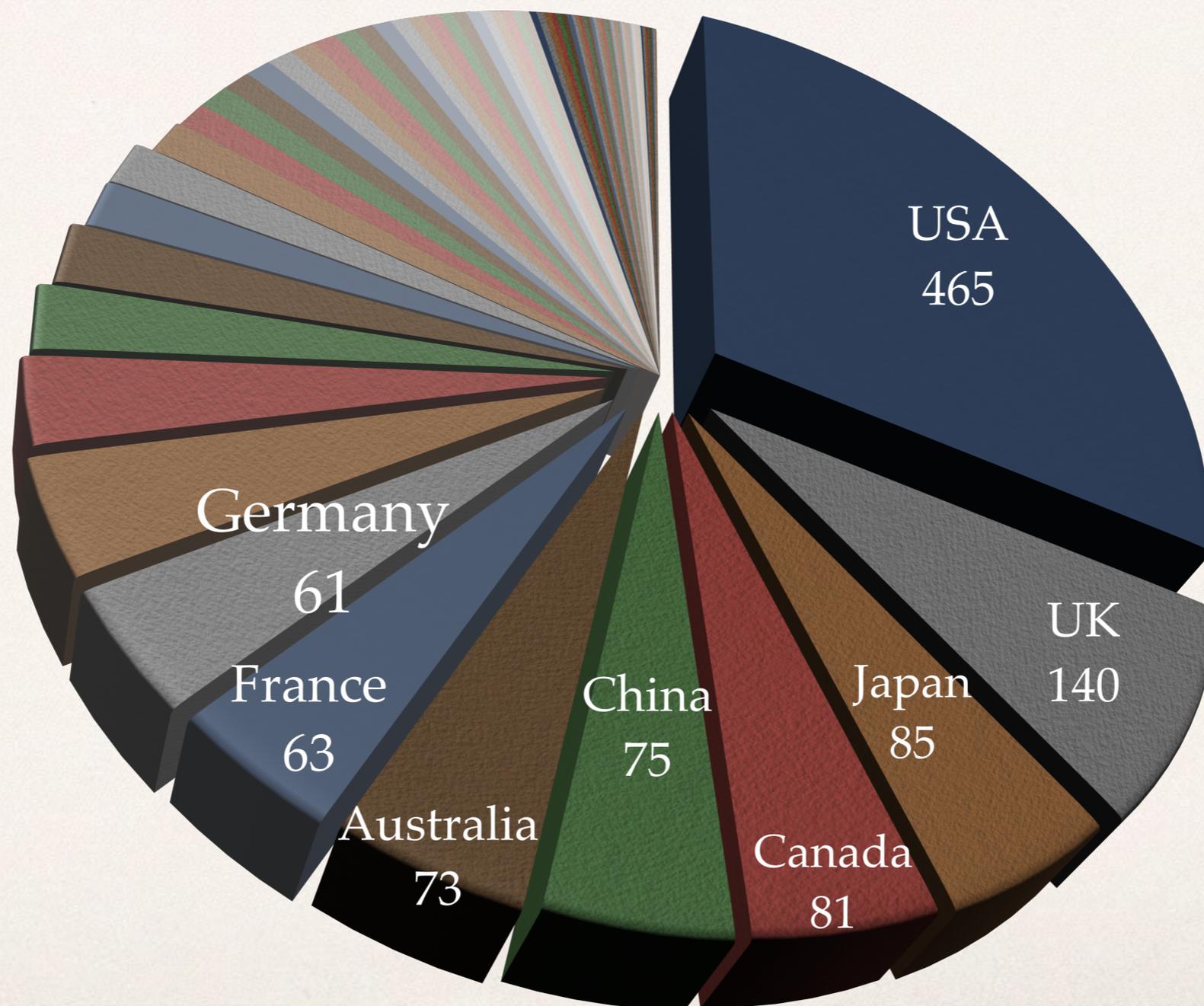
- ❖ **Phase 1:** Consultation with corpus linguists & civil engineers on the design & make-up of SCCERA (balanced & representative)
- ❖ **Phase 2:** Construction of SCCERA
- ❖ **Phase 3:** Quantitative & qualitative analysis of the corpus
- ❖ **Phase 4:** Exploring pedagogic applications of the corpus
- ❖ **Phase 5:** Dissemination of research results

SCCERA characteristics

- ❖ Total size: ~ 8 million words
- ❖ 45 journals (43 cited in SCI Expanded or SSCI)
- ❖ 1,100 research articles (most cited or downloaded)
- ❖ Year of publication: Range = 1989 - 2014; Mean = 2009
- ❖ 3,807 contributing authors (average of 3.46 authors per article)
- ❖ 1,598 participating institutions from 80 countries

Participating institutions by country (N = 80)

Participating institutions by country (N = 80)



Research questions

- ❖ What are the most frequently occurring keywords in civil engineering RAs?
- ❖ How are keywords from SCCERA distributed across established external wordlists - General Service List (West, 1953) & Academic Word List (Coxhead, 2000)?
- ❖ How are keywords from SCCERA distributed across the sub-divisions of civil engineering?
- ❖ What are the most frequent 3 - 6-word bundles in civil engineering RAs?
- ❖ How can this data inform materials design in an ESP context?

Keywords analysis (BNC reference corpus)

- * Highlights words which are unusually frequent / infrequent compared to a reference corpus & helps to characterize the 'aboutness' of a genre
- * Potentially more useful for materials design than raw word frequency data
- * SCCERA keywords focus on materials; observing, measuring & reporting on their physical properties or behaviours
- * Large degree of nominalisation (analysis, behaviour, distribution, etc), shifting the focus from an agent to an object or concept

Top 50 keywords in SCCERA

1	et al	26	temperature
2	fig	27	measured
3	model	28	behavior
4	data	29	coefficient
5	equation	30	ratio
6	results	31	spatial
7	values	32	variables
8	models	33	distribution
9	flow	34	strain
10	concrete	35	method
11	table	36	parameter
12	shear	37	measurements
13	using	38	shown
14	wave	39	earthquake
15	figure	40	value
16	surface	41	density
17	parameters	42	average
18	water	43	respectively
19	analysis	44	precipitation
20	eq	45	displacement
21	soil	46	effects
22	based	47	cement
23	stress	48	climate
24	observed	49	maximum
25	velocity	50	project

Top 40 keywords in SCCERA

et. al	table	soil	ratio
fig	shear	based	spatial
model	using	stress	variables
data	wave	the	distribution
equation	figure	observed	strain
results	surface	velocity	method
values	parameters	temperature	parameter
models	water	measured	measurements
flow	analysis	behavior	shown
concrete	eq	coefficient	earthquake

RANGE (Heatley & Nation, 1994) - Vocabulary profiling tool

	<u>Families</u>	<u>Types</u>	<u>Tokens</u>	<u>Percent</u>
K1 Words (1-1000):	99	134	142	29.40%
Function:	(13)	(2.69%)
Content:	(129)	(26.71%)
> Anglo-Sax =Not Greco-Lat/Fr Cog:	(37)	(7.66%)
K2 Words (1001-2000):	41	53	53	10.97%
> Anglo-Sax:	(20)	(4.14%)
1k+2k			...	(40.37%)
AWL Words (academic):	97	145	145	30.02%
> Anglo-Sax:	(4)	(0.83%)
Off-List Words:	<u>?</u>	<u>143</u>	<u>143</u>	<u>29.61%</u>
	237+?	475	483	100%

Current profile	
%	Cumul.
29.40	29.40
10.97	40.37
30.02	70.39
29.61	100.00

Distribution of keyword families across GSL, AWL & 'off-list'

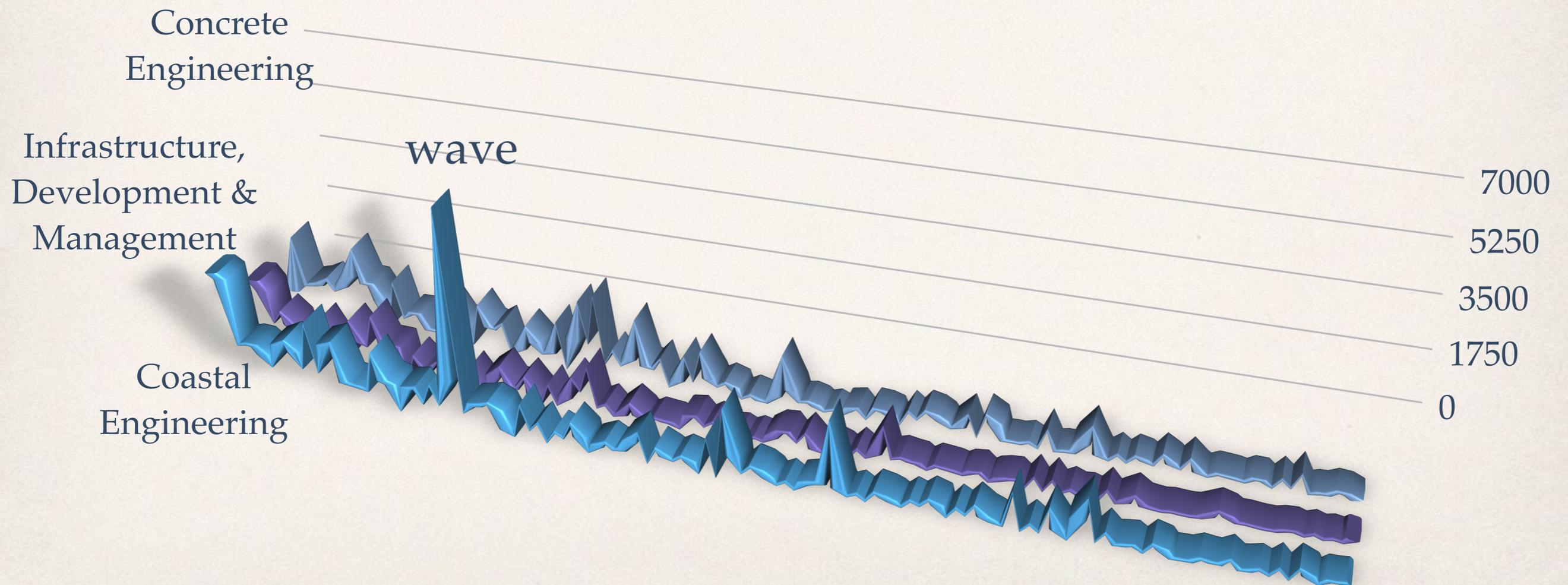


Distribution of keyword families across GSL, AWL & 'off-list'

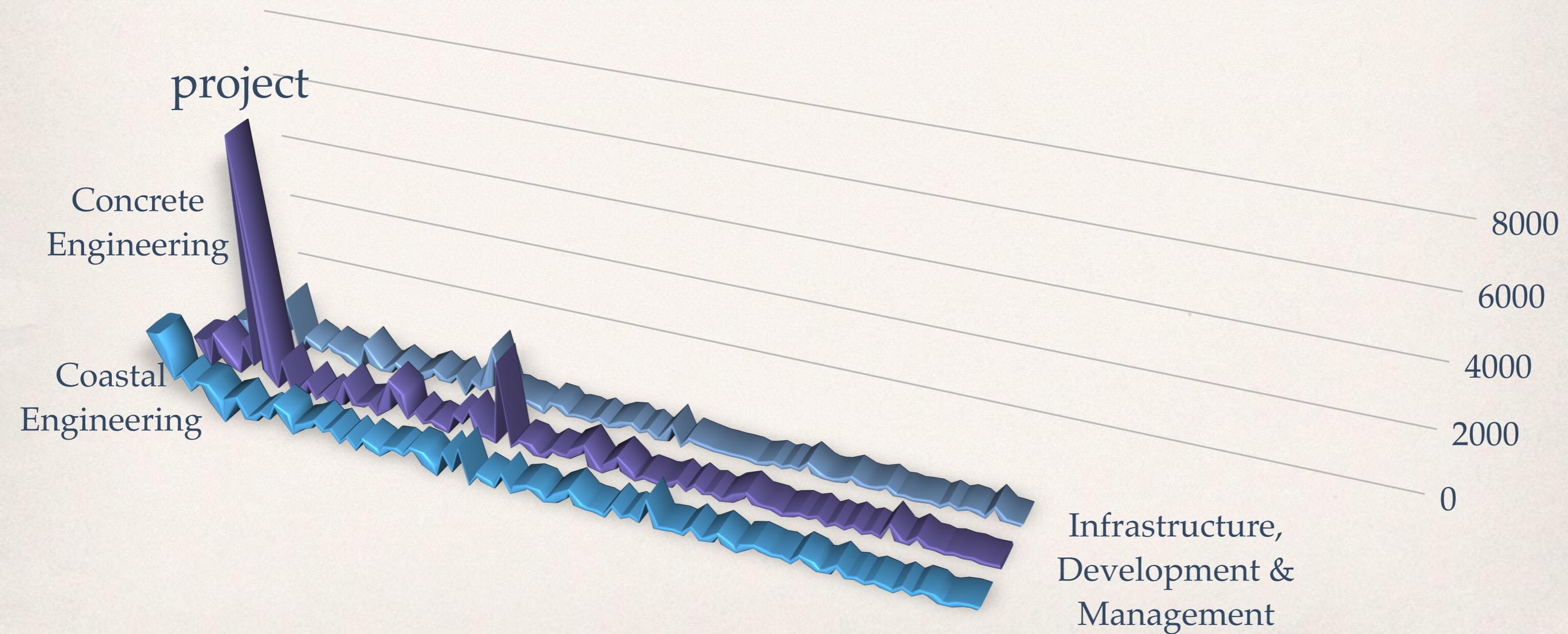
“Taken together, the first 2,000 words in West’s (1953) GSL and the word families in the AWL account for approximately 86% of the Academic Corpus” (Coxhead 2000: 222)



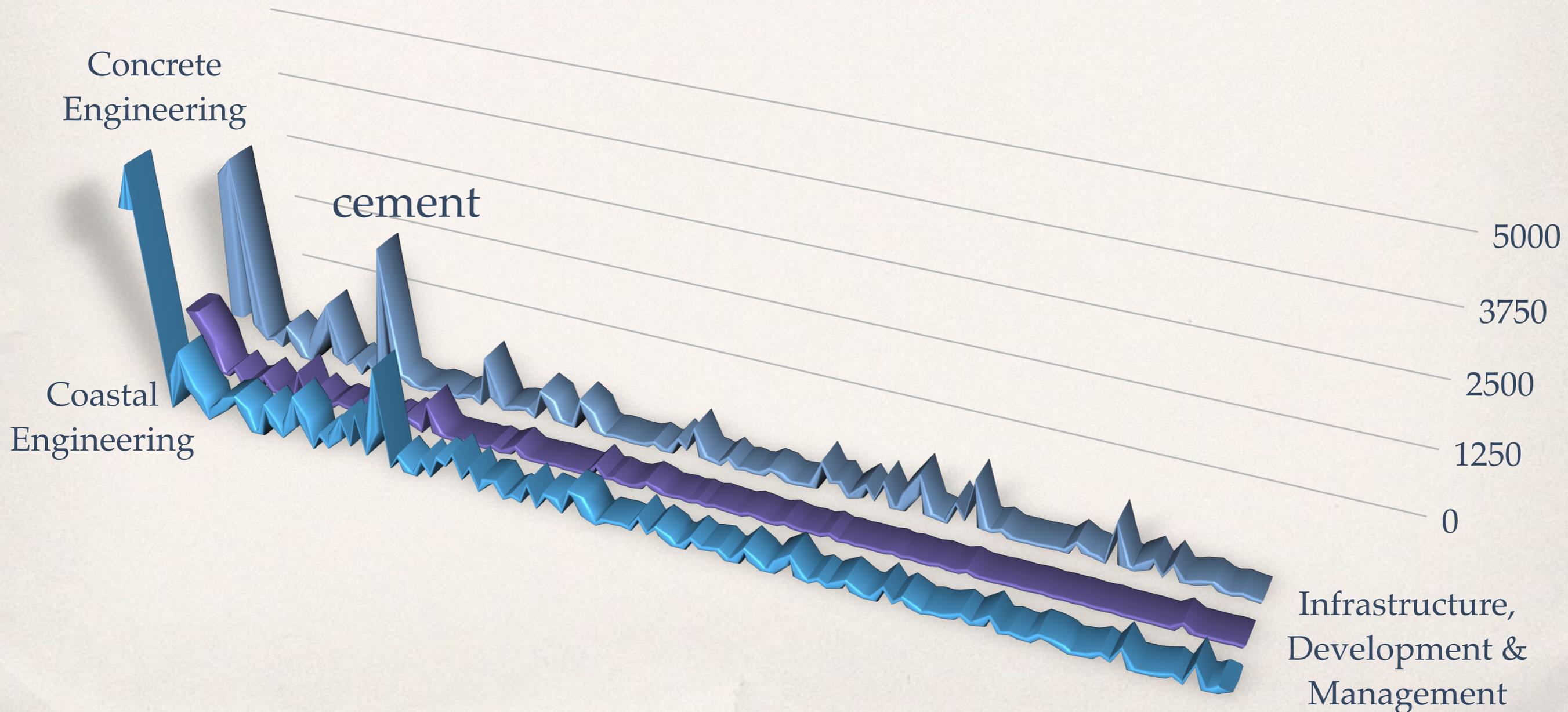
Distribution of keyword families across SCCERA: GSL



Distribution of keyword families across SCCERA: AWL



Distribution of keyword families across SCCERA: Offlist



AntConc - Wave collocates

AntConc 3.4.2m (Macintosh OS X) 2014

Corpus Files

- CCC_1 copy.txt
- CCC_2 copy.txt
- CCC_3 copy.txt
- CCC_4 copy.txt
- CCC_5 copy.txt
- CCC_6 copy.txt
- CCC_7 copy.txt
- CCC_8 copy.txt
- CCC_9 copy.txt
- CCC_10 copy.txt
- CCC_11 copy.txt
- CCC_12 copy.txt
- CCC_13 copy.txt
- CCC_14 copy.txt
- CCC_15 copy.txt
- CCC_16 copy.txt
- CCC_17 copy.txt
- CCC_18 copy.txt
- CCC_19 copy.txt
- CCC_20 copy.txt
- CCC_21 copy.txt
- CCC_22 copy.txt
- CCC_23 copy.txt
- CCC_24 copy.txt
- CCC_25 copy.txt
- CCR_1 copy.txt
- CCR_2 copy.txt
- CCR_3 copy.txt
- CCR_4 copy.txt
- CCR_5 copy.txt
- CCR_6 copy.txt
- CCR_7 copy.txt
- CCR_8 copy.txt
- CCR_9 copy.txt
- CCR_10 copy.txt
- CCR_11 copy.txt
- CCR_12 copy.txt
- CCR_13 copy.txt
- CCR_14 copy.txt
- CCR_15 copy.txt
- CCR_16 copy.txt
- CCR_17 copy.txt
- CCR_18 copy.txt
- CCR_19 copy.txt
- CCR_20 copy.txt
- CCR_21 copy.txt
- CCR_22 copy.txt
- CCR_23 copy.txt
- CCR_24 copy.txt
- CCR_25 copy.txt
- CE_1 copy.txt
- CE_2 copy.txt
- CE_3 copy.txt

Concordance Hits 6851

Hit	KWIC	File
5612	he MSE to approximate adequately	JWPCOE_11 cc
5613	through the gaps. To examine the	JWPCOE_8 cop
5614	, d two years after repair 4.2	MS_19 copy.t
5615	coherence effects resulting from	EESD_14 copy
5616	, 1995; Arduin et al., 2011b].	JGRO_9 copy.
5617	how wavelength would affect the	JWPCOE_11 cc
5618	of the shoal submergence on the	JWPCOE_11 cc
5619	wall jet scour (2D case) and	JWPCOE_17 cc
5620	n the scour hole development. In	JWPCOE_17 cc
5621	ilibrium for wall jet scour	JWPCOE_17 cc
5622	nd has d 50 =0.2mm . Fig. 14.	JWPCOE_17 cc
5623	del more accurate. Fig. 20.	JWPCOE_17 cc
5624	For turbulent wall jet scour and	JWPCOE_17 cc
5625	Scour-Hole Shape Evolution Three	JWPCOE_6 cop
5626	s. Mathematical Model The	JWPCOE_6 cop
5627	all breakwater (sometimes called	JWPCOE_8 cop
5628	many ways. Floating structures,	JCR_22 copy.
5629	creasing FB efficiency in short	JWPCOE_25 cc
5630	dle the interactions between the	CE_11 copy.t
5631	ions for the interactions of the	CE_11 copy.t
5632	500. Fig. 10. The sketch of the	CE_11 copy.t
5633	und. Keywords Integrated model;	CE_11 copy.t
5634	merical model to investigate the	CE_11 copy.t
5635	egrated model (PORO-WSSI II) for	CE_11 copy.t
5636	proposed an integrated model for	CE_11 copy.t
5637	is, 2005). The phenomenon of the	CE_11 copy.t
5638	associated with seabed structure	JGRO_5 copy

Search Term Words Case Regex

Search Window Size 50

Start Stop Sort

Kwic Sort

Level 1 1R Level 2 2R Level 3 3R

Clone Results

Total No. 1100
Files Processed

AntConc - *Wave* collocates (80+)

Collocate	Hits
height	813
velocity	279
period	263
break	249
energy	174
model	148
condition	147
propagation	145
runup	115

Lexical bundles

- ❖ Defined as “the most frequently occurring lexical sequences in a register” (Biber, Conrad & Cortes 2004)
- ❖ “lexical bundles are crucially important for the construction of discourse in all university registers” (Biber 2006: 174)

Lexical bundles are...

- ❖ Extremely common
- ❖ Usually not idiomatic in meaning
- ❖ Usually not perceptually salient
- ❖ Usually not complete structural units, but instead tend to bridge 2 phrases (e.g. ~5% for academic writing)
- ❖ Function as 'discourse frames' for the expression of new information - they don't express new propositional meaning themselves

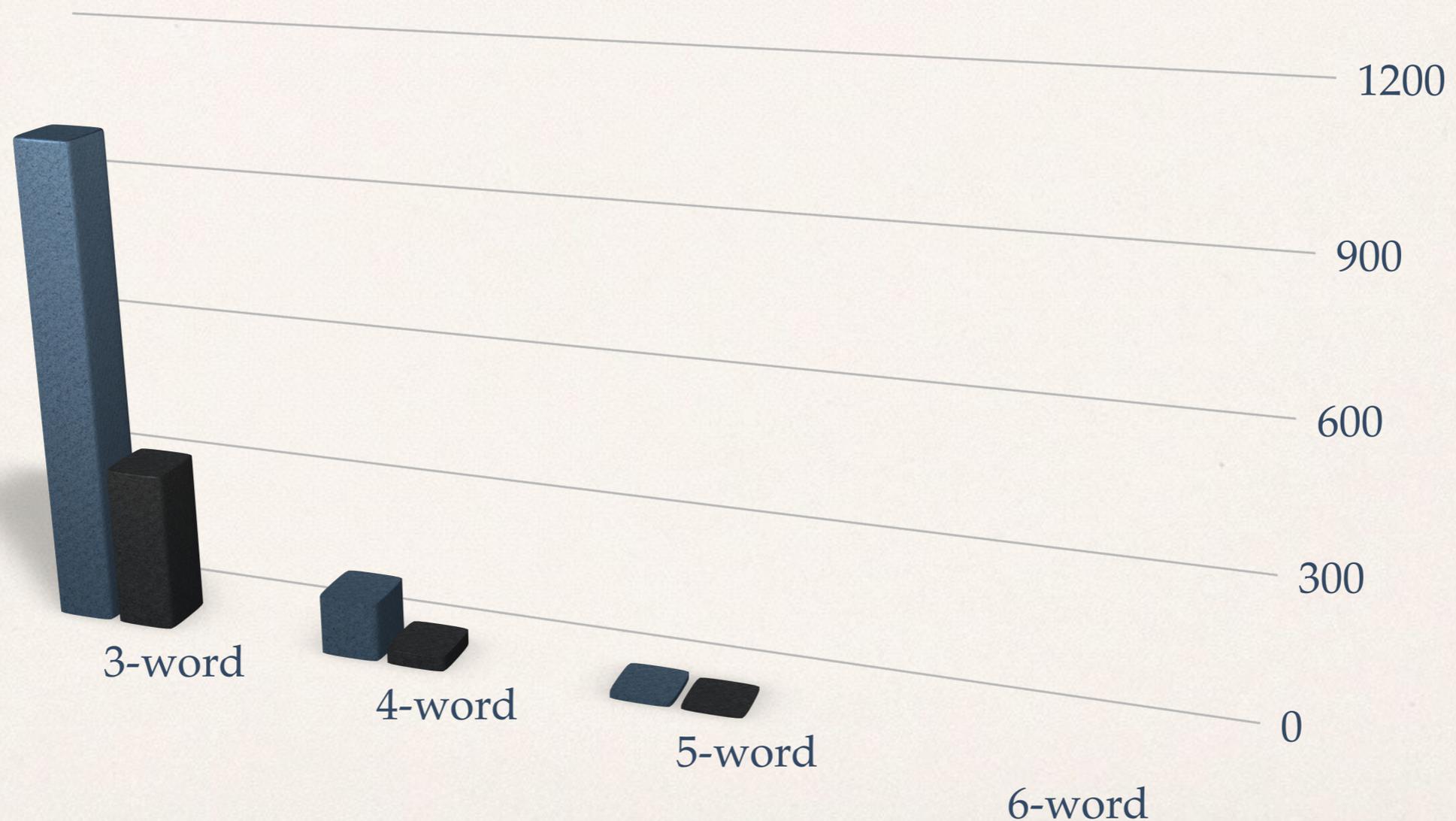
(Biber & Barbieri 2007: 269 / 70)

Lexical bundles (formulaic language)

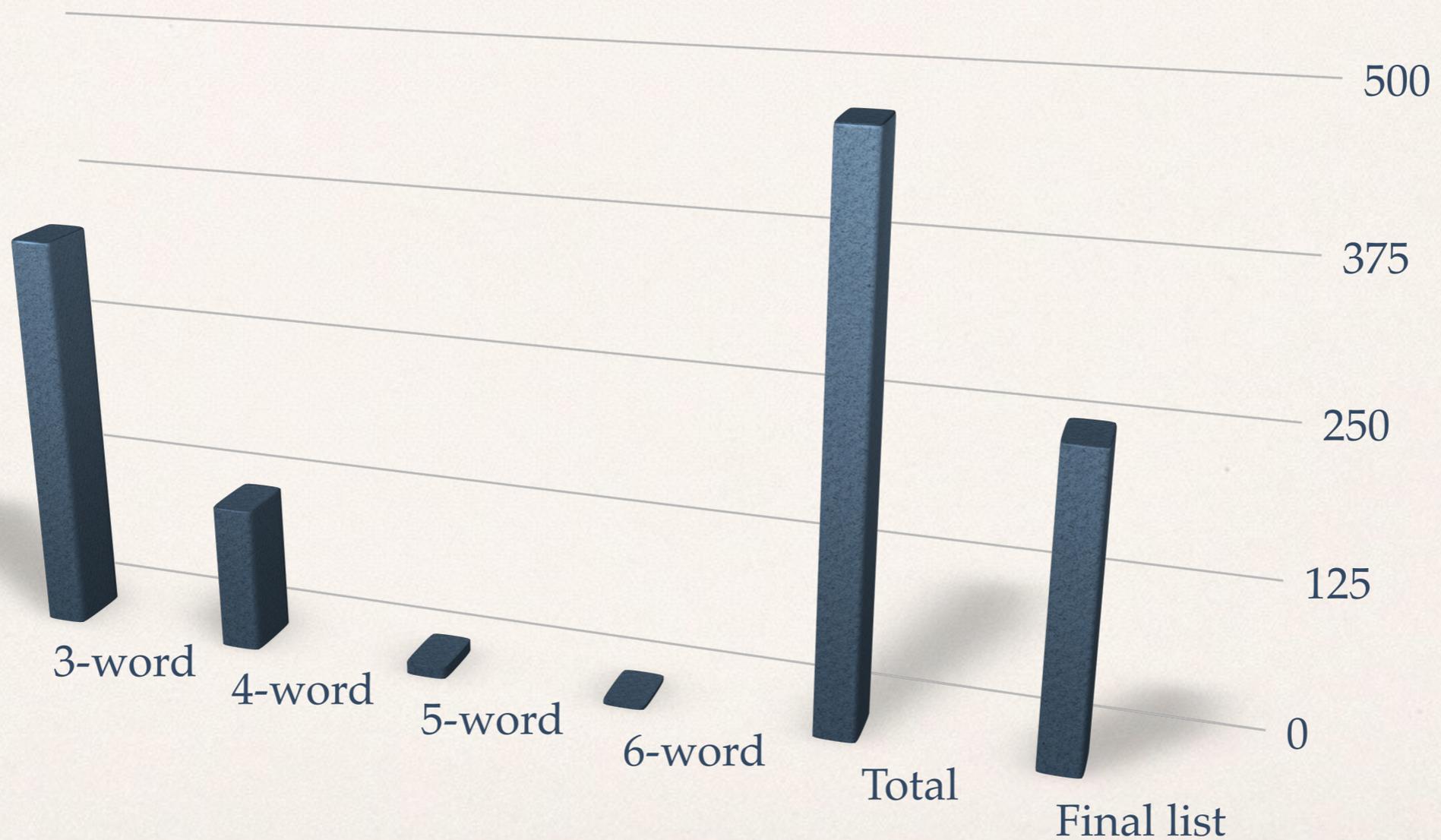
“Given the importance of formulaic language, it can be argued that it needs to be part of language syllabuses [...] Unfortunately, this is not the case. A perusal of almost any EFL/ESL textbook or test yields a paucity of formulaic sequences targeted for explicit attention/noticing, and even for those that do occur, there does not seem to be much principled basis for selection”

(Martinez & Schmitt 2012: 301)

Lexical bundles in SCCERA (20 or 40 occurrences per million cut-off)



Lexical bundles in SCCERA



Selecting lexical bundles for pedagogical applications

- ❖ Start with longest bundles (6-word), e.g. *it should be noted that the*
- ❖ Work through list from longest to shortest (6-word → 3-word), using key words (e.g. *noted*) to search for related bundles: *be noted that; should be noted*
- ❖ Use KWIC function in a concordancer (e.g. AntConc) to decide if bundles are fragments of longer bundles or not

Selecting lexical bundles - examples in the decision-making process

- * *should be noted* - excluded from list (KWIC search of SCCERA using AntConc indicates it is a fragment of *it should be noted that*)

should be noted: total hits = 374

Right sorted: *should be noted that* = 339 hits (90.6%)

Left sorted: *it should be noted* = 364 (97.3%)

- * *in the form* - included in final list (not solely a fragment of *in the form of*, but also *in the form + eqn*)

in the form: total hits = 395

Right sorted: *in the form of* = 335 hits (84.8%)

Selecting lexical bundles - examples in the decision-making process

* should indicate

should

Right s

Left sc

* in the f also in

in the f

Right s

The screenshot shows the AntConc 3.4.2m (Macintosh OS X) 2014 interface. On the left, a 'Corpus Files' list contains 24 files (CE_1.txt to CE_24.txt) and 12 files (JCR_1.txt to JCR_12.txt). The main window displays 'Concordance Hits 374' for the search term 'should be noted'. The results are shown in a table with columns 'Hit', 'KWIC', and 'File'. The KWIC column shows the search term highlighted in blue, with surrounding text in other colors. The search options at the bottom include 'Search Term' (should be noted), 'Words' checked, 'Case' unchecked, 'Regex' unchecked, 'Search Window Size' (50), and 'Kwic Sort' (Level 1 3R, Level 2 4R, Level 3 5R). A 'Clone Results' button is also visible.

Hit	KWIC	File
24	ver on a decadal basis. It should be noted, however, th	RSE_12.txt
25	relationships discussed.) It should be noted, however, th	JGGE_21.txt
26	screpancy is not known. It should be noted, however, th	JACS_2.txt
27	, limitations of MIP method should be noted in evaluatic	CCC_19.txt
28	its energy requirements. It should be noted in passing t	JEM_4.txt
29	chm. On the other hand, it should be noted in the propo	TRB_12.txt
30	differences are very slight should be noted. It means th	MS_17.txt
31	The causes of human injury should be noted not only by	JNDS_8.txt
32	e resulting paired readings should be noted on a field t	JCR_4.txt
33	t well pass unremarked. It should be noted that a conce	GT_12.txt
34	l distribution function. It should be noted that a direc	SF_18.txt
35	le quantities de and d~. It should be noted that a displ	JEM_5.txt
36	ydration times. However, it should be noted, that a dis	CCR_3.txt
37	for comparison purposes. It should be noted that a large	JGGE_17.txt
38	d raw data is obtained. It should be noted that a perfe	TGRS_20.txt
39	ent replacement levels. It should be noted that a plair	CCC_24.txt

AntConc

of, but

Selecting lexical bundles - examples in the decision-making process

- * *should be noted* - excluded from list (KWIC search of SCCERA using AntConc indicates it is a fragment of *it should be noted that*)

should be noted: total hits = 374

Right sorted: *should be noted that* = 339 hits (90.6%)

Left sorted: *it should be noted* = 364 (97.3%)

- * *in the form* - included in final list (not solely a fragment of *in the form of*, but also *in the form + eqn*)

in the form: total hits = 395

Right sorted: *in the form of* = 335 hits (84.8%)

Selecting lexical bundles - examples in the decision-making process

* should be indicated

should

Right side

Left side

* in the foreground also in

in the foreground

Right side

The screenshot shows the AntConc 3.4.2m (Macintosh OS X) 2014 interface. On the left, a 'Corpus Files' list contains 25 files: CE_1.txt to CE_25.txt and JCR_1.txt to JCR_12.txt. Below the list, it shows 'Total No. 1100' and 'Files Processed' with a progress bar. The main window displays 'Concordance Hits 374' for the search term 'should be noted'. The results are shown in a table with columns 'Hit', 'KWIC', and 'File'. The KWIC column shows the search term highlighted in blue within its context. The 'File' column lists the source file for each hit. At the bottom, search options are visible: 'Search Term' is 'should be noted', 'Words' is checked, 'Case' and 'Regex' are unchecked. 'Search Window Size' is set to 50. 'Kwic Sort' options are checked for Level 1 (1L), Level 2 (2L), and Level 3 (3L). Buttons for 'Start', 'Stop', 'Sort', and 'Clone Results' are present.

Hit	KWIC	File
1	e columns. ACKNOWLEDGMENTS should be noted that the A&S	JSE_3.txt
2	Fe-containing analogues. In should be noted that the rec	CCR_10.txt
3	The causes of human injury should be noted not only by	JNDS_8.txt
4	n in the lattice). Also, is should be noted that if diff	JACS_8.txt
5	ossible improvements of ISPH should be noted here. Ataie-	JHR_15.txt
6	S-H phases (Fig. A.6). It should be noted that C-S-H	CCR_3.txt
7	vitational acceleration. It should be noted that Eq. (1	CE_2.txt
8	e validity was achieved. It should be noted that not al	JCEM_22.txt
9	led concrete aggregates. It should be noted that the FRA	CCC_21.txt
10	cial volume of pore air. It should be noted that pore ai	SF_3.txt
11	96, Antoine et al.1997). It should be noted that these t	IJGIS_15.txt
12	-1 [Wada et al., 2010]. It should be noted that the sim	WRR_7.txt
13	in Cetin et al. (2000). It should be noted that the pro	JGGE_21.txt
14	L-Qorchi et al. (2003)). It should be noted, however, th	JDE_32.txt
15	ss of run-up algorithm. It should be noted that a wave	CE_8.txt
16	nd selection algorithms. It should be noted that climat	RSE_13.txt

AntConc

of, but

Selecting lexical bundles - examples in the decision-making process

- * *should be noted* - excluded from list (KWIC search of SCCERA using AntConc indicates it is a fragment of *it should be noted that*)

should be noted: total hits = 374

Right sorted: *should be noted that* = 339 hits (90.6%)

Left sorted: *it should be noted* = 364 (97.3%)

- * *in the form* - included in final list (not solely a fragment of *in the form of*, but also *in the form + eqn*)

in the form: total hits = 395

Right sorted: *in the form of* = 335 hits (84.8%)

Selecting lexical bundles - examples in the decision-making process

* *should be indicate*

should be

Right so

Left so

* *in the for also in t*

in the fo

Right so

AntConc 3.4.2m (Macintosh OS X) 2014

Corpus Files

- CE_1.txt
- CE_2.txt
- CE_3.txt
- CE_4.txt
- CE_5.txt
- CE_6.txt
- CE_7.txt
- CE_8.txt
- CE_9.txt
- CE_10.txt
- CE_11.txt
- CE_12.txt
- CE_13.txt
- CE_14.txt
- CE_15.txt
- CE_16.txt
- CE_17.txt
- CE_18.txt
- CE_19.txt
- CE_20.txt
- CE_21.txt
- CE_22.txt
- CE_23.txt
- CE_24.txt
- CE_25.txt
- JCR_1.txt
- JCR_2.txt
- JCR_3.txt
- JCR_4.txt
- JCR_5.txt
- JCR_6.txt
- JCR_7.txt
- JCR_8.txt
- JCR_9.txt
- JCR_10.txt
- JCR_11.txt
- JCR_12.txt

Total No. 1100
Files Processed

Concordance Hits 395

Hit	KWIC	File
374	s of coastal tsunami impact in the form of video digital	JWPCOE_13.tx
375	best is to be of order 11 in the form equation(60)	JSV_6.txt
376	, G*(w,T)G*(w,T), in the form equation(73) w	JSV_6.txt
377	coefficients are presented in the form Equation (7) whe	NHR_16.txt
378	nd Jangid and Banerji [387] in the form equation(68) w	JSV_6.txt
379	small deflection was given in the form equation(49) w	JSV_6.txt
380	prediction step is written in the form Equation (5) whe	JHE_9.txt
381	ed well by a power function in the form: Equation (1) wr	SF_25.txt
382	ll rate model was developed in the form Equation (1) whe	NHR_16.txt
383	red the free-energy density in the form Equation (65) wr	JEM_12.txt
384	žant 1984) can be written in the form Equation (17) wr	CGJ_23.txt
385	curing time is approximated in the form. Equation (3) wr	SF_17.txt
386	nal expression may be given in the form: Equation (1) w	SCHM_14.txt
387	explicit scheme, is written in the form Equation (4) whe	JHE_9.txt
388	om Eq. 14 , repeated here in the form, Equation (16) w	JGGE_20.txt
389	cum equation can be written in the form Equation (1) whe	JHE_9.txt

Search Term Words Case Regex

Search Window Size 50

Start Stop Sort

Kwic Sort Level 1 4R Level 2 5R Level 3 6R

Clone Results

AntConc

of, but

Selecting lexical bundles - examples in the decision-making process

- * *should be noted* - excluded from list (KWIC search of SCCERA using AntConc indicates it is a fragment of *it should be noted that*)

should be noted: total hits = 374

Right sorted: *should be noted that* = 339 hits (90.6%)

Left sorted: *it should be noted* = 364 (97.3%)

- * *in the form* - included in final list (not solely a fragment of *in the form of*, but also *in the form + eqn*)

in the form: total hits = 395

Right sorted: *in the form of* = 335 hits (84.8%)

Selecting lexical bundles - examples in the decision-making process

- * *should be noted* - excluded from list (KWIC search of SCCERA using AntConc

indi

For plasticity with isotropic softening, Borino et al. (1999) considered the free-energy density in the form

sho

$$\rho\psi(\boldsymbol{\epsilon}, \boldsymbol{\epsilon}_p, \bar{\kappa}) = \rho\psi_e(\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_p) + \rho\psi_p(\bar{\kappa}) \quad (65)$$

Rig

where $\rho\psi_e(\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_p) = \frac{1}{2}(\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_p)^T \mathbf{D}_e(\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_p)$ = elastically stored energy and $\rho\psi_p(\bar{\kappa})$ = plastic part of the free energy, usually interpreted as the energy stored in microstructural changes. The postulate of maximum plastic dissipation then leads to an associated nonlocal plasticity model with the softening law (38) replaced by

Lef

- * *in the form of* also

in the

$$\sigma_Y = \sigma_0 + \widetilde{h(\bar{\kappa})} \quad (66)$$

Right sorted: *in the form of* = 335 hits (84.8%)

Selecting lexical bundles - examples in the decision-making process

- * *should be noted* - excluded from list (KWIC search of SCCERA using AntConc indicates it is a fragment of *it should be noted that*)

should be noted: total hits = 374

Right sorted: *should be noted that* = 339 hits (90.6%)

Left sorted: *it should be noted* = 364 (97.3%)

- * *in the form* - included in final list (not solely a fragment of *in the form of*, but also *in the form + eqn*)

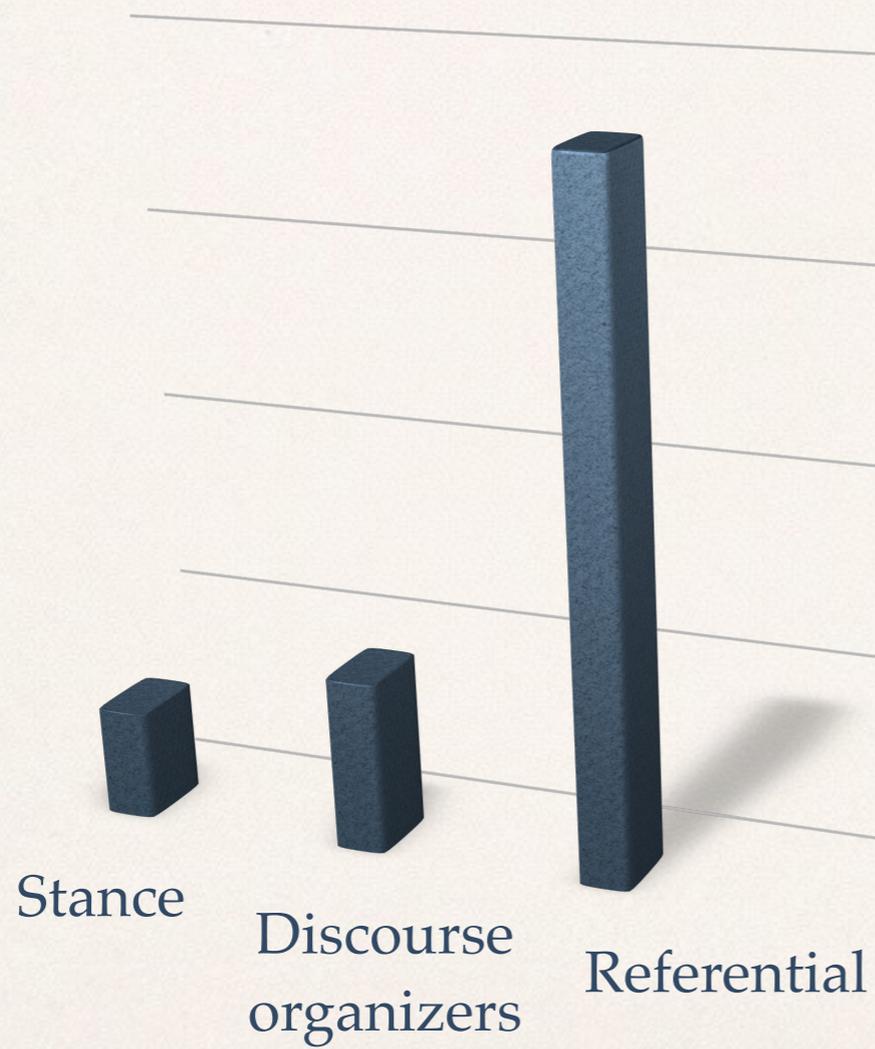
in the form: total hits = 395

Right sorted: *in the form of* = 335 hits (84.8%)

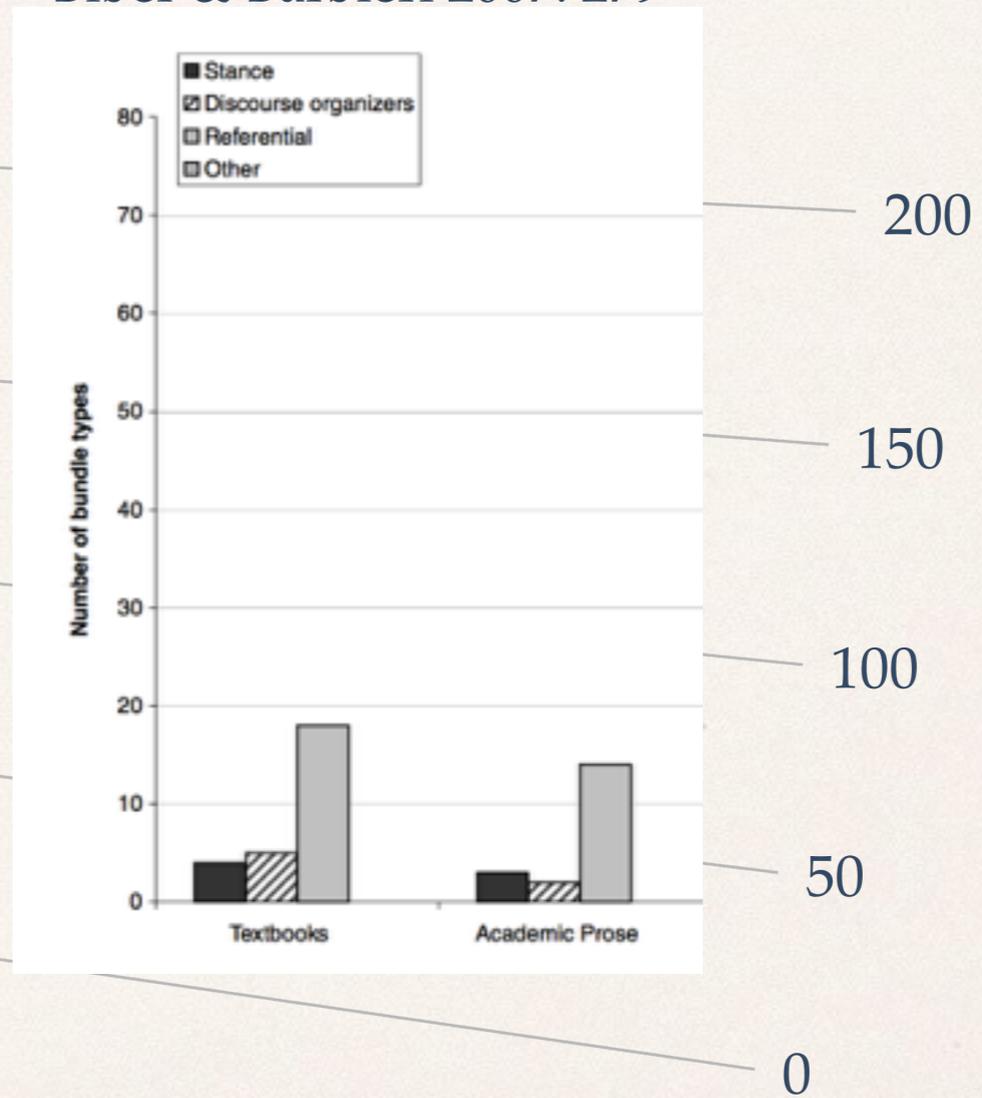
Typical patterns in clusters

- ❖ “Almost 30 per cent of the top 100 chains (types) have this pattern *PREP the N of the*. If we count closely related chains, with an indefinite article (*in the context of a*) or with embedded adjectives, such as *the ADJ N of the (the far side of the)*, then, over 40 per cent of the top 100 chains have this pattern.” (Stubbs & Barth 2003: 82)
- ❖ (34.7% in SCCERA)

Lexical bundles in SCCERA by type



Biber & Barbieri 2007: 279



Referential expressions

- ❖ Specifying attributes: *a little bit of; the size of the; the nature of the*
- ❖ Time / place / text deixis: *in the USA; at the end of the; as shown in Fig*

Conclusion

- * Keywords from SCCERA (compared with GSL & AWL) provide the foundation for a more systematic approach to academic vocabulary development for civil engineering students
- * Lexical bundles complement keyword lists in materials design by providing a 'phraseological profile' of a genre (Romer 2010: 27)
- * Because lexical bundles are frequent but not perceptually salient, 'they might be good candidates for overt instruction' (Biber & Barbieri 2007: 284)
- * "it is necessary to go beyond frequencies and complement the corpus-derived information with research tapping into their pedagogic relevance, with suggestions from teachers and students" (Flowerdew 2012: 268)

References

- Biber, D. (2006). *University Language. A corpus-based study of spoken and written registers*. Amsterdam: John Benjamins
- Biber, D. & Barbieri, F. (2007). 'Lexical bundles in university spoken and written registers'. *English for Specific Purposes* 26: 263-286
- Biber, D., Conrad, S. & Cortes, V. (2004). 'If you look at...: Lexical bundles in university teaching and textbooks'. *Applied Linguistics* 25.3: 371-405
- Coxhead, A. (2000). 'A new academic word list'. *TESOL Quarterly* 34.2: 213-238.
- Flowerdew, L. (2012). *Corpora and language education*. Basingstoke: Palgrave Macmillan
- Heatley, A. & P. Nation (1994). 'Range', Victoria University of Wellington, NZ.
- Hyland, K. (2008). 'As can be seen: Lexical bundles and disciplinary variation'. *English for Specific Purposes* 27: 4-21
- Hyland, K. & Milton, J. (1997). 'Qualification and certainty in L1 and L2 students' writing'. *Journal of Second Language Writing* 6(2): 183-205

References

- Martinez, R. & N. Schmitt (2012). 'A phrasal expressions list'. *Applied Linguistics* 33.3: 299-320
- Pearson, J. (1998). *Terms in context*. Amsterdam: John Benjamins
- Rea Rizzo, C. (2010). 'Getting on with corpus compilation: From theory to practice'. *ESP World* 1(9): 1-22
- Römer, U. (2010). 'Using general and specialised corpora in English language teaching: Past, present and future'. In Campoy-Cubillo, Belles-Fortunato & Gea-Valor (Eds.) *Corpus-based approaches to English language teaching*. London: Continuum: 18-35
- Scott, M. (2011). *WordSmith Tools Version 6*, Liverpool: Lexical Analysis Software.
- Stubbs, M. & I. Barth (2003). 'Using recurrent phrases as text-type discriminators'. *Functions of Language* 10.1: 61-104
- West, M. (1953). *A general service list of English words*. London: Longman, Green & Co.







Lexical bundles

- ❖ Lexical bundles occurring at 20-40 times per million words or more considered (conservative - less conservative approaches: see Biber & Barbieri 2007)
- ❖ “an ESP perspective, considering each register on its own terms, is required to adequately describe the use of lexical bundles in the university context” (Biber & Barbieri 2007: 265)
- ❖ Three main functions of lexical bundles (Biber & Barbieri 2007: 270):
 1. Expressions of stance (epistemic, desire, obligation, intention/prediction, ability)
 2. Discourse organizers (topic introduction, topic elaboration/clarification, referential identification/focus)
 3. Referential expressions (imprecision indicators, specifying attributes, time/place/text-deixis) - dominant function in written university registers

Conclusion

- ❖ Keywords from SCCERA (compared with GSL & AWL) provide the foundation for a systematic approach to academic vocabulary development for civil engineering students
- ❖ Word bundles???
- ❖ “A student in the university context will frequently encounter [word bundles]... failure to understand their textual and interpersonal functions will obviously influence her/his success in dealing with both spoken and written language situations” (Biber & Barbieri 2007: 284)
- ❖ “more work probably has to be put into the creation of reliable corpus-based language descriptions for learners and teachers, especially descriptions of specialized discourses” (Römer 2010: 29)

Conclusion

- ❖ “[...] some researchers have argued that factors such as perceptual salience and developmental readiness are more important than frequency for acquisition... If these findings are generalizable, they might indicate that the discourse functions of lexical bundles are not easily acquired: they are not perceptually salient, despite their high frequency. If so, lexical bundles might be good candidates for overt instruction.” (Biber & Barbieri 2007: 284)
- ❖ “it would seem sensible to combine our formulaic sequences into these wordlists in a way that would create a much more inclusive overall description of the most frequent (and therefore useful) lexical items of English, both individual and multi-word” (Martinez & Schmitt 2012: 304)

Lexical bundles can be deceptively transparent

- ❖ “It is not unreasonable to guess that L2 learners processing those expressions [in the PHRASE List] might therefore actually believe they understand them (if they identify them) simply because the individual words are so well known, making them, in Laufer’s (1989) terms, ‘deceptively transparent’ (Martinez & Schmitt 2012: 313)

Lexical bundles included in assessments of complexity?

- ❖ “There is also the potential [...] for the development of an automated lexical profiling tool such as *Range* (Heatley and Nation 1994) that instead of only analyzing a text for individual words [...] also carries out a ‘sweep’ for phrases [...] to more accurately reflect its lexical complexity’ (Martinez & Schmitt 2012: 315/6)

Lexical bundles in SCCERA compared to PHRASE List

- ❖ Only **13.2%** of lexical bundles in SCCERA occur in Martinez & Schmitt's (2012) PHRASal Expressions List (PHRASE List) - 505 most frequent, non-transparent multi-word expressions in English (derived from BNC)
- ❖ Specialized corpus (SCCERA) vs general corpus (BNC)
- ❖ Transparency (receptive vs productive goals)
- ❖ 3- to 6-word bundles vs 2- to 4-word bundles

3-word clusters

based on the	with respect to	in the case
as well as	in this paper	there is a
the number of	one of the	the value of
in order to	in this study	the presence of
shown in fig	a function of	can be used
in terms of	the case of	the fact that
due to the	part of the	according to the
the effect of	a number of	as a result
the use of	the effects of	be used to
as shown in	the results of	the other hand

4-word clusters

in the case of	the results of the	it is important to
on the other hand	is shown in fig	it should be noted
as a function of	the size of the	in the context of
as shown in fig	are shown in fig	is assumed to be
as well as the	is based on the	the fact that the
can be used to	the end of the	should be noted that
on the basis of	at the end of	in the form of
with respect to the	the effect of the	it is possible to
in terms of the	at the same time	it can be seen
as a result of	in the united states	in this paper we

SCCERA journal list

Journal	Article Title	Article Code	Number Authors	Institution Countries	Year of Publication	Number Words
Coastal Engineering	Modelling storm impacts on beaches, dunes	CE_1	6	Netherlands; USA	2009	11,398
Coastal Engineering	Corrected Incompressible SPH method for a	CE_2	3	Japan; UK	2008	8,184
Coastal Engineering	44-year wave hindcast for the North East At	CE_3	3	Portugal; Spain	2008	2,817
Coastal Engineering	Increasing wave heights and extreme value	CE_4	3	USA	2010	10,302
Coastal Engineering	Modified Moving Particle Semi-implicit meth	CE_5	1	Japan	2009	11,452
Coastal Engineering	Numerical analysis of wave overtopping of r	CE_6	4	Spain	2008	7,208
Coastal Engineering	A 44-year high-resolution ocean and atmos	CE_7	5	Spain; France	2008	8,829
Coastal Engineering	Simulation of nonlinear wave run-up with a	CE_8	2	Denmark	2008	7,499
Coastal Engineering	Beach Wizard: Nearshore bathymetry estim	CE_9	6	Netherlands; USA; Chile	2008	8,286
Coastal Engineering	Efficient computation of surf zone waves usi	CE_10	2	Netherlands	2008	6,532
Coastal Engineering	An integrated model for the wave-induced s	CE_11	4	UK; China; USA	2013	8,876
Coastal Engineering	Statistical simulation of wave climate and e	CE_12	4	Australia	2008	8,888
Coastal Engineering	Hindcast of the wave conditions along the w	CE_13	3	Portugal	2008	5,480
Coastal Engineering	Laboratory and numerical studies of wave d	CE_14	3	USA	2009	6,662
Coastal Engineering	A probabilistic methodology to estimate futu	CE_15	3	UK	2008	6,178
Coastal Engineering	Run-up of tsunamis and long waves in term	CE_16	2	Denmark	2008	8,813
Coastal Engineering	Measurement of wave-by-wave bed-levels in	CE_17	3	Australia; UK	2008	1,648
Coastal Engineering	The morphological response of a nearshore	CE_18	4	Netherlands; USA	2008	7,938
Coastal Engineering	Direct bed shear stress measurements in bc	CE_19	4	Australia; UK	2009	9,281
Coastal Engineering	Two-dimensional time dependent hurricane	CE_20	7	Netherlands; USA	2010	9,300
Coastal Engineering	Modeling hurricane waves and storm surge	CE_21	10	Netherlands; USA	2011	9,528
Coastal Engineering	On the evolution and run-up of breaking sol	CE_22	4	Taiwan	2008	7,007
Coastal Engineering	Morphodynamic responses to the deep water	CE_23	3	China	2009	7,962
Coastal Engineering	Large-scale dune erosion tests to study the	CE_24	5	Netherlands	2008	5,852
Coastal Engineering	Wave boundary layer over a stone-covered	CE_25	4	Denmark	2008	11,155
J. of Coastal Research	The Role of Remote Sensing in Predicting ar	JCR_1	1	USA	2009	7,568
J. of Coastal Research	Shoreline Definition and Detection: A Review	JCR_2	2	Australia	2005	5,895
J. of Coastal Research	Erosion Hazard Vulnerability of US Coastal C	JCR_3	3	USA	2005	5,147
J. of Coastal Research	A Simple Method of Measuring Beach Profile	JCR_4	2	Portugal	2006	2,389
J. of Coastal Research	A New Global Coastal Database for Impact a	JCR_5	9	Greece; UK; Ireland; Germany; N	2008	4,788
J. of Coastal Research	Assessment of Vulnerability and Adaptation	JCR_6	1	Germany	2008	8,688
J. of Coastal Research	Sustainable Management of Surfing Breaks:	JCR_7	4	New Zealand	2009	11,534
J. of Coastal Research	Importance of Coastal Change Variables in I	JCR_8	3	USA	2010	4,311
J. of Coastal Research	The Healing Sea: A Sustainable Coastal Oce	JCR_9	2	Belgium	2009	11,155
J. of Coastal Research	Open-Ocean Barrier Islands: Global Influen	JCR_10	2	USA	2011	7,538
J. of Coastal Research	Tracking Oil Slicks and Predicting their Traje	JCR_11	1	USA	2010	6,583
J. of Coastal Research	Classification of Coasts	JCR_12	1	USA	2004	8,476
J. of Coastal Research	Coastal Classification: Systematic Approach	JCR_13	1	USA	2004	20,515

Phase 1: Consultation on corpus design

Corpus linguists & academics from 12 departments of Civil Engineering consulted on design criteria:

- * Peer-reviewed journals, preferably listed in Science Citation Index Expanded (SCI) or Social Sciences Citation Index (SSCI)
- * Widely read & respected by researchers; considered “key journals” or “desired outlets for academic work”
- * Higher impact factors (IF), 5-year IF, Eigenfactor, article influence (Thomson Reuters)
- * Research articles selected by (a) Most cited; (b) Most viewed; (c) Most recent (1 article per volume)
- * Minimum size of 1 million words recommended for specialised corpora (Kennedy 1998; Pearson 1998; Rea Rizzo 2010)

Phase 2: Construction of SCCERA

- ❖ HTML or PDF version of articles copied into MS Word
- ❖ Extraneous information removed (references, date of acceptance, author affiliation, contact info., tables & figures, equations)
- ❖ Text cleaned up using spelling & grammar checking function of MS Word (hyphenated words, conjoined words, character misreadings)
- ❖ HTML fragments ('Table options', 'Turn Mathjax on', etc.) removed using find & replace function in MS Word
- ❖ Articles saved as text-only (.txt) files

Phase 2: Construction of SCCERA

- ❖ 2nd round of cleaning up using text-only files (Greek symbols, etc.)
- ❖ Final document checked against original PDF file
- ❖ SCCERA part-of-speech (POS) tagged using CLAWS 4 (Lancaster University UCREL C7 tag set (Total no. tag types = 137):
<http://ucrel.lancs.ac.uk/claws7tags.html>)

Phase 2: Construction of SCCERA

Coastal Engineering 56 (2009) 1133–1152

Contents lists available at ScienceDirect

Coastal Engineering

journal homepage: www.elsevier.com/locate/coastaleng

Modelling storm impacts on beaches, dunes and barrier islands

Dano Roelvink ^{a,h,c,*}, Ad Reniers ^{c,d}, Ap van Dongeren ^b, Jaap van Thiel de Vries ^{b,c}, Robert McCall ^{b,c}, Jamie Lescinski ^b

^a UNESCO-IHE Institute for Water Education, P.O. Box 3015, 2601 DA Delft, The Netherlands
^b Delft, The Netherlands
^c Delft University of Technology, The Netherlands
^d Rosenstiel School of Marine and Atmospheric Science, Univ. of Miami, United States

ARTICLE INFO

ABSTRACT

Article history:
Received 15 December 2008
Received in revised form 12 July 2009
Accepted 18 August 2009
Available online 15 September 2009

Keywords:
Swash
Low-frequency waves
Dune erosion
Overtopping
Overwashing
Breaching
Morphology

A new nearshore numerical model approach to assess the natural coastal response during time-varying storm and hurricane conditions, including dune erosion, overwash and breaching, is validated with a series of analytical, laboratory and field test cases. Innovations include a non-stationary wave driver with directional spreading to account for wave-group generated surf and swash motions and an avalanching mechanism providing a smooth and robust solution for slumping of sand during dune erosion. The model performs well in different situations including dune erosion, overwash and breaching with specific emphasis on swash dynamics, avalanching and 2DH effects; these situations are all modelled using a standard set of parameter settings. The results show the importance of infragravity waves in extending the reach of the resolved processes to the dune front. The simple approach to account for slumping of the dune face by avalanching makes the model easily applicable in two dimensions and applying the same settings good results are obtained both for dune erosion and breaching.

© 2009 Elsevier B.V. All rights reserved.



```
<text>_NULL

^ Modelling_VVG storm_NN1 impacts_NN2 on_II beaches_NN2 ,,, dunes_NN2 and_CC
barrier_NN1 islands_NN2 Abstract_VV@ A_ZZ1 new_JJ nearshore_NN1 numerical_JJ
model_NN1 approach_NN1 to_TO assess_VVI the_AT natural_JJ coastal_JJ
response_NN1 during_II time-varying_JJ storm_NN1 and_CC hurricane_NN1
conditions_NN2 ,,, including_II dune_NN1 erosion_NN1 ,,, overwash_NN1 and_CC
breaching_VVG ,,, is_VBZ validated_VVN@ with_IW a_AT1 series_NN of_IO
analytical_JJ ,,, laboratory_NN1 and_CC field_NN1 test_NN1 cases_NN2 ._.

^ Innovations_NN2 include_VV@ a_AT1 non-stationary_JJ wave_NN1 driver_NN1
with_IW directional_JJ spreading_NN1 to_TO account_VVI for_IF wave-group_JJ
generated_JJ@ surf_NN1 and_CC swash_VVI motions_NN2 and_CC an_AT1
avalanching_JJ@ mechanism_NN1 providing_VVG a_AT1 smooth_JJ and_CC robust_JJ
solution_NN1 for_IF slumping_VVG of_IO sand_NN1 during_II dune_NN1 erosion_NN1
'..'

^ The_AT model_NN1 performs_VVZ well_RR in_II different_JJ situations_NN2
including_II dune_NN1 erosion_NN1 ,,, overwash_NN1 and_CC breaching_VVG
with_IW specific_JJ emphasis_NN1 on_II swash_NN1 dynamics_NN ,,,
avalanching_VVG and_CC 2DH_FO effects_NN2 ;;; these_DD2 situations_NN2 are_VBR
all_DB modelled_VVN using_VVG a_AT1 standard_JJ set_NN1 of_IO parameter_NN1
settings_NN2 ._.

^ The_AT results_NN2 show_VV@ the_AT importance_NN1 of_IO infragravity_NN1
waves_NN2 in_II extending_VVG the_AT reach_NN1 of_IO the_AT resolved_JJ@
processes_NN2 to_II the_AT dune_NN1 front_NN1 ._.

^ The_AT simple_JJ approach_NN1 to_TO account_VVI for_IF slumping_VVG of_IO
the_AT dune_NN1 face_NN1 by_II avalanching_VVG makes_VVZ the_AT model_NN1
easily_RR applicable_JJ in_II two_MC dimensions_NN2 and_CC applying_VVG the_AT
same_DA settings_NN2 good_JJ results_NN2 are_VBR obtained_VVN both_RR for_IF
dune_NN1 erosion_NN1 and_CC breaching_VVG ._.

```

UCREL CLAWS 7 Tagset

VVG = -ing participle of lexical verb

NN1 = singular common noun

NN2 = plural common noun

II = general preposition

Processing time (mins per RA)

→ Mean = 7.5

March 10
2014

April 19
2014

Processing time (mins per RA)



Common problems

1. Introduction

A primary goal of modeling physical processes in the atmospheric and hydrologic sciences is the prediction of a variable in time and/or space from a given set of inputs. How well a model fits the observed data (referred to as model evaluation, or sometimes as model validation) usually is determined by pairwise comparisons of model-simulated (or model-predicted) values with observations. Quantitative assessments of the degree to which the model simulations match the observations are used to provide an evaluation of the model's predictive abilities.

Frequently, evaluations of model performance utilize a number of statistics and techniques. Usually included in these tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient (r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997].

As knowledge of physical processes has increased, models have become more complex. Often these models include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.

Paper number 1998WR900018.
0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled until the differences between the observed and model-simulated data are minimized [Nash and Sutcliffe, 1970; Song and James, 1991; Hay, 1998]. Stochastic calibration procedures are usually employed, which limits graphical analyses of scatterplots, for example, so that statistical analyses must be solely used. Consequently, statistics other than r and R^2 have been developed to describe better the degree of association between the observed and model-simulated data. The objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

In this paper, three basic methods for model evaluation will be discussed: the coefficient of determination R^2 , the coefficient of efficiency E [Nash and Sutcliffe, 1970], and the index of agreement d [Willmott et al., 1985]. In general, this paper addresses comparisons of model-simulated data (P) with the observed data (O) for the same set of conditions (i.e., a pairwise comparison) over a given time period divided into N time increments that can be of arbitrary duration (e.g., monthly or daily time steps).

2.1. Coefficient of Determination R^2

The coefficient of determination is the square of the Pearson's product-moment correlation coefficient (i.e., $R^2 = r^2$) and describes the proportion of the total variance in the observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

1. Introduction

A primary goal of modeling physical processes in the atmospheric and hydrologic sciences is the prediction of a variable in time and/or space from a given set of inputs. How well a model fits the observed data (referred to as model evaluation, or sometimes as model validation) usually is determined by pairwise comparisons of model-simulated (or model-predicted) values with observations. Quantitative assessments of the degree to which the model simulations match the observations are used to provide an evaluation of the model's predictive abilities.

Frequently, evaluations of model performance utilize a number of statistics and techniques. Usually included in these tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient

(r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of

a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997].

As knowledge of physical processes has increased, models have become more complex. Often these models

include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.

Paper number 1998WR900018. 0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled

until the differences between the observed and model-

simulated data are minimized [Nash and Sutcliffe, 1970; Song

and James, 1991; Hay, 1998]. Stochastic calibration procedures

are usually employed, which limits graphical analyses of scatter-

plots, for example, so that statistical analyses must be solely

used. Consequently, statistics other than r and R^2 have been

developed to describe better the degree of association between the observed and model-simulated data. The

objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated

with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

Common problems

Words
split with
hyphens

1. Introduction

A primary goal of modeling physical processes in the atmospheric and hydrologic sciences is the prediction of a variable in time and/or space from a given set of inputs. How well a model fits the observed data (referred to as model evaluation, or sometimes as model validation) usually is determined by pairwise comparisons of model-simulated (or model-predicted) values with observations. Quantitative assessments of the degree to which the model simulations match the observations are used to provide an evaluation of the model's predictive abilities.

Frequently, evaluations of model performance utilize a number of statistics and techniques. Usually included in these tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient (r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997].

As knowledge of physical processes has increased, models have become more complex. Often these models include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.

Paper number 1998WR900018.
0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled until the differences between the observed and model-simulated data are minimized [Nash and Sutcliffe, 1970; Song and James, 1991; Hay, 1998]. Stochastic calibration procedures are usually employed, which limits graphical analyses of scatterplots, for example, so that statistical analyses must be solely used. Consequently, statistics other than r and R^2 have been developed to describe better the degree of association between the observed and model-simulated data. The objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

In this paper, three basic methods for model evaluation will be discussed: the coefficient of determination R^2 , the coefficient of efficiency E [Nash and Sutcliffe, 1970], and the index of agreement d [Willmott et al., 1985]. In general, this paper addresses comparisons of model-simulated data (P) with the observed data (O) for the same set of conditions (i.e., a pairwise comparison) over a given time period divided into N time increments that can be of arbitrary duration (e.g., monthly or daily time steps).

2.1. Coefficient of Determination R^2

The coefficient of determination is the square of the Pearson's product-moment correlation coefficient (i.e., $R^2 = r^2$) and describes the proportion of the total variance in the observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

1. Introduction

A primary goal of modeling physical processes in the atmospheric and hydrologic sciences is the prediction of a variable in time and/or space from a given set of inputs. How well a model fits the observed data (referred to as model evaluation, or sometimes as model validation) usually is determined by pairwise comparisons of model-simulated (or model-predicted) values with observations. Quantitative assessments of the degree to which the model simulations match the observations are used to provide an evaluation of the model's predictive abilities. Frequently, evaluations of model performance utilize a number of statistics and techniques. Usually included in these tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient

(r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of

a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997].

As knowledge of physical processes has increased, models have become more complex. Often these models

include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.

Paper number 1998WR900018. 0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled

until the differences between the observed and model-

simulated data are minimized [Nash and Sutcliffe, 1970; Song

and James, 1991; Hay, 1998]. Stochastic calibration procedures

are usually employed, which limits graphical analyses of scatter-

plots, for example, so that statistical analyses must be solely

used. Consequently, statistics other than r and R^2 have been

developed to describe better the degree of association between the observed and model-simulated data. The

objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated

with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

Common problems

Words split with hyphens

Specialised words incorrectly identified as mistakes (collinearity)

1. Introduction

A primary goal of modeling physical processes in the atmospheric and hydrologic sciences is the prediction of a variable in time and/or space from a given set of inputs. How well a model fits the observed data (referred to as model evaluation, or sometimes as model validation) usually is determined by pairwise comparisons of model-simulated (or model-predicted) values with observations. Quantitative assessments of the degree to which the model simulations match the observations are used to provide an evaluation of the model's predictive abilities.

Frequently, evaluations of model performance utilize a number of statistics and techniques. Usually included in these tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient (r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997].

As knowledge of physical processes has increased, models have become more complex. Often these models include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.

Paper number 1998WR900018.
0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled until the differences between the observed and model-simulated data are minimized [Nash and Sutcliffe, 1970; Song and James, 1991; Hay, 1998]. Stochastic calibration procedures are usually employed, which limits graphical analyses of scatterplots, for example, so that statistical analyses must be solely used. Consequently, statistics other than r and R^2 have been developed to describe better the degree of association between the observed and model-simulated data. The objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

In this paper, three basic methods for model evaluation will be discussed: the coefficient of determination R^2 , the coefficient of efficiency E [Nash and Sutcliffe, 1970], and the index of agreement d [Willmott et al., 1985]. In general, this paper addresses comparisons of model-simulated data (P) with the observed data (O) for the same set of conditions (i.e., a pairwise comparison) over a given time period divided into N time increments that can be of arbitrary duration (e.g., monthly or daily time steps).

2.1. Coefficient of Determination R^2

The coefficient of determination is the square of the Pearson's product-moment correlation coefficient (i.e., $R^2 = r^2$) and describes the proportion of the total variance in the observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

1. Introduction

A primary goal of modeling physical processes in the variable in time and/or space from a given set of input model evaluation, or sometimes as model validation) simulated (or model-predicted) values with observations model simulations match the observations are used to Frequently, evaluations of model performance utilize number of statistics and techniques. Usually included tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient (r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of

a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997]. As knowledge of physical processes has increased, models have become more complex. Often these models include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.
Paper number 1998WR900018. 0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled until the differences between the observed and model-simulated data are minimized [Nash and Sutcliffe, 1970; Song and James, 1991; Hay, 1998]. Stochastic calibration procedures are usually employed, which limits graphical analyses of scatterplots, for example, so that statistical analyses must be solely

used. Consequently, statistics other than r and R have been developed to describe better the degree of association between the observed and model-simulated data. The objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

Common problems

Words split with hyphens

Specialised words incorrectly identified as mistakes (collinearity)

Letters incorrectly identified (ll)

1. Introduction

A primary goal of modeling physical processes in the atmospheric and hydrologic sciences is the prediction of a variable in time and/or space from a given set of inputs. How well a model fits the observed data (referred to as model evaluation, or sometimes as model validation) usually is determined by pairwise comparisons of model-simulated (or model-predicted) values with observations. Quantitative assessments of the degree to which the model simulations match the observations are used to provide an evaluation of the model's predictive abilities.

Frequently, evaluations of model performance utilize a number of statistics and techniques. Usually included in these tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient (r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997].

As knowledge of physical processes has increased, models have become more complex. Often these models include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.

Paper number 1998WR900018.
0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled until the differences between the observed and model-simulated data are minimized [Nash and Sutcliffe, 1970; Song and James, 1991; Hay, 1998]. Stochastic calibration procedures are usually employed, which limits graphical analyses of scatterplots, for example, so that statistical analyses must be solely used. Consequently, statistics other than r and R^2 have been developed to describe better the degree of association between the observed and model-simulated data. The objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

In this paper, three basic methods for model evaluation will be discussed: the coefficient of determination R^2 , the coefficient of efficiency E [Nash and Sutcliffe, 1970], and the index of agreement d [Willmott et al., 1985]. In general, this paper addresses comparisons of model-simulated data (P) with the observed data (O) for the same set of conditions (i.e., a pairwise comparison) over a given time period divided into N time increments that can be of arbitrary duration (e.g., monthly or daily time steps).

2.1. Coefficient of Determination R^2

The coefficient of determination is the square of the Pearson's product-moment correlation coefficient (i.e., $R^2 = r^2$) and describes the proportion of the total variance in the observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

1. Introduction

A primary goal of modeling physical processes in the variable in time and/or space from a given set of input model evaluation, or sometimes as model validation) simulated (or model-predicted) values with observations model simulations match the observations are used to Frequently, evaluations of model performance utilize number of statistics and techniques. Usually included tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient (r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of

a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997]. As knowledge of physical processes has increased, models have become more complex. Often these models include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.
Paper number 1998WR900018. 0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled until the differences between the observed and model-simulated data are minimized [Nash and Sutcliffe, 1970; Song and James, 1991; Hay, 1998]. Stochastic calibration procedures are usually employed, which limits graphical analyses of scatterplots, for example, so that statistical analyses must be solely used. Consequently, statistics other than r and R have been developed to describe better the degree of association between the observed and model-simulated data. The objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

Common problems

1. Introduction

A primary goal of modeling physical processes in the atmospheric and hydrologic sciences is the prediction of a variable in time and/or space from a given set of inputs. How well a model fits the observed data (referred to as model evaluation, or sometimes as model validation) usually is determined by pairwise comparisons of model-simulated (or model-predicted) values with observations. Quantitative assessments of the degree to which the model simulations match the observations are used to provide an evaluation of the model's predictive abilities.

Frequently, evaluations of model performance utilize a number of statistics and techniques. Usually included in these tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient (r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997].

As knowledge of physical processes has increased, models have become more complex. Often these models include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.

Paper number 1998WR900018.
0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled until the differences between the observed and model-simulated data are minimized [Nash and Sutcliffe, 1970; Song and James, 1991; Hay, 1998]. Stochastic calibration procedures are usually employed, which limits graphical analyses of scatterplots, for example, so that statistical analyses must be solely used. Consequently, statistics other than r and R^2 have been developed to describe better the degree of association between the observed and model-simulated data. The objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

In this paper, three basic methods for model evaluation will be discussed: the coefficient of determination R^2 , the coefficient of efficiency E [Nash and Sutcliffe, 1970], and the index of agreement d [Willmott et al., 1985]. In general, this paper addresses comparisons of model-simulated data (P) with the observed data (O) for the same set of conditions (i.e., a pairwise comparison) over a given time period divided into N time increments that can be of arbitrary duration (e.g., monthly or daily time steps).

2.1. Coefficient of Determination R^2

The coefficient of determination is the square of the Pearson's product-moment correlation coefficient (i.e., $R^2 = r^2$) and describes the proportion of the total variance in the observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

1. Introduction

A primary goal of modeling physical processes in the variable in time and/or space from a given set of input model evaluation, or sometimes as model validation) simulated (or model-predicted) values with observations model simulations match the observations are used to Frequently, evaluations of model performance utilize number of statistics and techniques. Usually included tools are "goodness-of-fit" or relative error measures (bounded statistics, usually between 0.0 and 1.0) to assess the ability of a model to simulate reality. Often these statistics are based on the familiar Pearson's product-moment correlation coefficient (r) or its square, the coefficient of determination (R^2). These two statistics describe the degree of collinearity between the observed and model-simulated variates. They are almost always discussed in basic statistics texts and, consequently, are familiar to virtually all scientists. Unfortunately, both r and R^2 suffer from limitations that make them poor measures of model performance. Although these statistics continue to be used to determine how well a model simulates the observed data, they nevertheless provide a biased view of the efficacy of

a model [Willmott, 1981; Willmott et al., 1985; Kessler and Neas, 1994; Legates and Davis, 1997]. As knowledge of physical processes has increased, models have become more complex. Often these models include numerous parameters that are calibrated through optimization

Copyright 1999 by the American Geophysical Union.

Paper number 1998WR900018. 0043-1397/99/1998WR900018\$09.00

procedures, where a range in model parameters is sampled until the differences between the observed and model-simulated data are minimized [Nash and Sutcliffe, 1970; Song and James, 1991; Hay, 1998]. Stochastic calibration procedures are usually employed, which limits graphical analyses of scatterplots, for example, so that statistical analyses must be solely used. Consequently, statistics other than r and R^2 have been developed to describe better the degree of association between the observed and model-simulated data. The objectives of this paper are to (1) examine various goodness-of-fit measures and to identify limitations associated with each, and (2) suggest viable alternative measures for the evaluation of hydrologic and hydroclimatic models.

2. Statistics for Evaluation of Hydrologic and Hydroclimatic Models

Words split with hyphens

Specialised words incorrectly identified as mistakes (collinearity)

Letters incorrectly identified (ll)

Text broken up by footnotes, page numbers, etc.

Common problems

Copyright 1999 by the American Geophysical Union.
 Paper number 1998WR900018.
 0043-1397/99/1998WR900018\$09.00

observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

233

234

LEGATES AND MCCABE: EVALUATING "GOODNESS-OF-FIT" MEASURES

$$R^2 = \frac{\left[\sum_{i=1}^N (O_i - \bar{O})(P_i - \bar{P}) \right]^2}{\left[\sum_{i=1}^N (O_i - \bar{O})^2 \right] \left[\sum_{i=1}^N (P_i - \bar{P})^2 \right]} \quad (1)$$

where the overbar denotes the mean for the entire time period of the evaluation. Note, however, that the coefficient of determination is limited in that it standardizes for differences between the observed and predicted means and variances since it

adjusting factor would result in an increase in the correlation, possibly causing it to exceed 1.0 in extreme cases. Consequently, we do not advocate the use of such adjusting factors.

It should be noted that nonparametric or rank correlation methods also exist (e.g., Spearman's rho or Kendall's tau). As nonparametric statistics, they are less sensitive to outliers in the data and generally provide a more robust characterization of the correlation between observed and predicted values. Unfortunately, rank correlation measures are associated with a loss of information as interval/ratio data are converted to ordinal (ranked) form [see *Burt and Barber, 1996*], and, like their parametric counterparts, they are not sensitive to additive and proportional differences between the observed and model-simulated values.

and describes the proportion of the total variance in the observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

233

)

234

LEGATES AND MCCABE: EVALUATING "GOODNESS-OF-FIT" MEASURES

2 R=

$$2 \frac{\sum_{i=1}^N (O_i - \bar{O})(P_i - \bar{P})}{\left[\sum_{i=1}^N (O_i - \bar{O})^2 \right] \left[\sum_{i=1}^N (P_i - \bar{P})^2 \right]}$$

N

1

adjusting factor would result in an increase in the correlation, possibly causing it to exceed 1.0 in extreme cases. Consequently, we do not advocate the use of such adjusting factors.

It should be noted that nonparametric or rank correlation methods also exist (e.g., Spearman's rho or Kendall's tau). As nonparametric statistics, they are less sensitive to outliers in the data and generally provide a more robust characterization of the correlation between observed and predicted values. Unfortunately, rank correlation measures are associated with a loss of information as interval/ratio data are converted to ordinal (ranked) form [see *Burt and Barber, 1996*], and, like their parametric counterparts, they are not sensitive to additive and proportional differences between the observed and model-simulated values.

2.2. Coefficient of Efficiency E

The coefficient of efficiency *E* has been widely used to evaluate the performance of hydrologic models [e.g., *Leavesley et al, 1983; Wdcox et al, 1990*]. *Nash and Sutcliffe* [1970] defined the coefficient of efficiency which ranges from minus infinity to 1.0, with higher values indicating better agreement, as

[...] and is given by (Equation 1)
 where x denotes...

Common problems

Copyright 1999 by the American Geophysical Union.
 Paper number 1998WR900018.
 0043-1397/99/1998WR900018\$09.00

observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

233

Mathematical symbols
not recognised

234

LEGATES AND MCCABE: EVALUATING "GOODNESS-OF-FIT" MEASURES

$$R^2 = \frac{\left[\sum_{i=1}^N (O_i - \bar{O})(P_i - \bar{P}) \right]^2}{\left[\sum_{i=1}^N (O_i - \bar{O})^2 \right] \left[\sum_{i=1}^N (P_i - \bar{P})^2 \right]} \quad (1)$$

where the overbar denotes the mean for the entire time period of the evaluation. Note, however, that the coefficient of determination is limited in that it standardizes for differences between the observed and predicted means and variances since it

adjusting factor would result in an increase in the correlation, possibly causing it to exceed 1.0 in extreme cases. Consequently, we do not advocate the use of such adjusting factors.

It should be noted that nonparametric or rank correlation methods also exist (e.g., Spearman's rho or Kendall's tau). As nonparametric statistics, they are less sensitive to outliers in the data and generally provide a more robust characterization of the correlation between observed and predicted values. Unfortunately, rank correlation measures are associated with a loss of information as interval/ratio data are converted to ordinal (ranked) form [see *Burt and Barber, 1996*], and, like their parametric counterparts, they are not sensitive to additive and proportional differences between the observed and model-simulated values.

and describes the proportion of the total variance in the observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

233

)

234

LEGATES AND MCCABE: EVALUATING "GOODNESS-OF-FIT" MEASURES

2 R=

$$2 \frac{\sum_{i=1}^N (O_i - \bar{O})(P_i - \bar{P})}{\left[\sum_{i=1}^N (O_i - \bar{O})^2 \right] \left[\sum_{i=1}^N (P_i - \bar{P})^2 \right]}$$

N

1

adjusting factor would result in an increase in the correlation, possibly causing it to exceed 1.0 in extreme cases. Consequently, we do not advocate the use of such adjusting factors.

It should be noted that nonparametric or rank correlation methods also exist (e.g., Spearman's rho or Kendall's tau). As nonparametric statistics, they are less sensitive to outliers in the data and generally provide a more robust characterization of the correlation between observed and predicted values. Unfortunately, rank correlation measures are associated with a loss of information as interval/ratio data are converted to ordinal (ranked) form [see *Burt and Barber, 1996*], and, like their parametric counterparts, they are not sensitive to additive and proportional differences between the observed and model-simulated values.

2.2. Coefficient of Efficiency E

The coefficient of efficiency *E* has been widely used to evaluate the performance of hydrologic models [e.g., *Leavesley et al, 1983; Wdcox et al, 1990*]. *Nash and Sutcliffe* [1970] defined the coefficient of efficiency which ranges from minus infinity to 1.0, with higher values indicating better agreement, as

[...] and is given by (Equation 1)
where x denotes...

Common problems

Copyright 1999 by the American Geophysical Union.
 Paper number 1998WR900018.
 0043-1397/99/1998WR900018\$09.00

observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

233

234

LEGATES AND MCCABE: EVALUATING "GOODNESS-OF-FIT" MEASURES

$$R^2 = \frac{\left[\sum_{i=1}^N (O_i - \bar{O})(P_i - \bar{P}) \right]^2}{\left[\sum_{i=1}^N (O_i - \bar{O})^2 \right] \left[\sum_{i=1}^N (P_i - \bar{P})^2 \right]} \quad (1)$$

where the overbar denotes the mean for the entire time period of the evaluation. Note, however, that the coefficient of determination is limited in that it standardizes for differences between the observed and predicted means and variances since it

adjusting factor would result in an increase in the correlation, possibly causing it to exceed 1.0 in extreme cases. Consequently, we do not advocate the use of such adjusting factors.

It should be noted that nonparametric or rank correlation methods also exist (e.g., Spearman's rho or Kendall's tau). As nonparametric statistics, they are less sensitive to outliers in the data and generally provide a more robust characterization of the correlation between observed and predicted values. Unfortunately, rank correlation measures are associated with a loss of information as interval/ratio data are converted to ordinal (ranked) form [see *Burt and Barber*, 1996], and, like their parametric counterparts, they are not sensitive to additive and proportional differences between the observed and model-simulated values.

and describes the proportion of the total variance in the observed data that can be explained by the model. It ranges from 0.0 to 1.0, with higher values indicating better agreement, and is given by

233

234

LEGATES AND MCCABE: EVALUATING "GOODNESS-OF-FIT" MEASURES

234

$$R^2 = \frac{\left[\sum_{i=1}^N (O_i - \bar{O})(P_i - \bar{P}) \right]^2}{\left[\sum_{i=1}^N (O_i - \bar{O})^2 \right] \left[\sum_{i=1}^N (P_i - \bar{P})^2 \right]}$$

adjusting factor would result in an increase in the correlation, possibly causing it to exceed 1.0 in extreme cases. Consequently, we do not advocate the use of such adjusting factors.

It should be noted that nonparametric or rank correlation methods also exist (e.g., Spearman's rho or Kendall's tau). As nonparametric statistics, they are less sensitive to outliers in the data and generally provide a more robust characterization of the correlation between observed and predicted values. Unfortunately, rank correlation measures are associated with a loss of information as interval/ratio data are converted to ordinal (ranked) form [see *Burt and Barber*, 1996], and, like their parametric counterparts, they are not sensitive to additive and proportional differences between the observed and model-simulated values.

2.2. Coefficient of Efficiency *E*

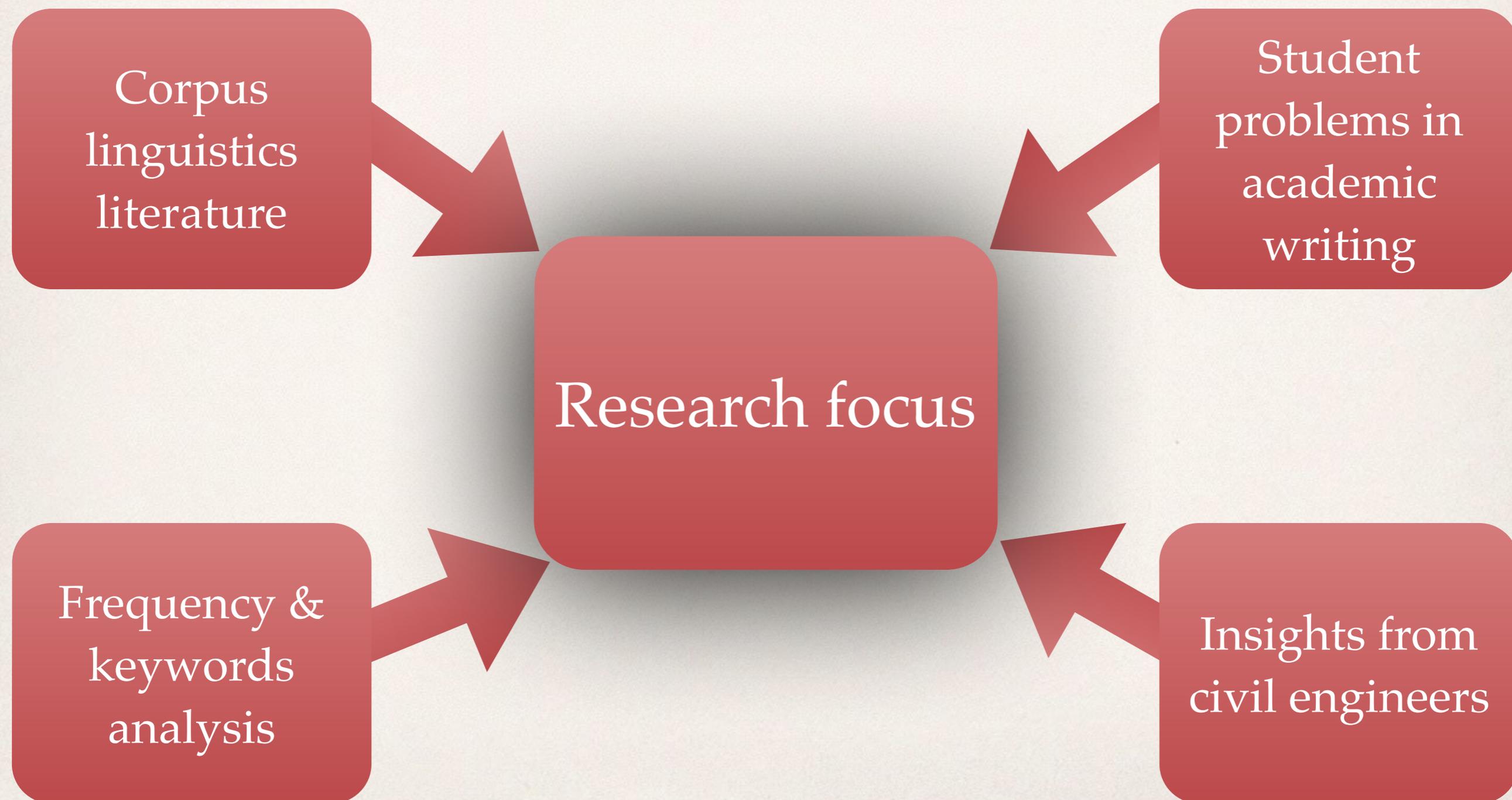
The coefficient of efficiency *E* has been widely used to evaluate the performance of hydrologic models [e.g., *Leavesley et al*, 1983; *Wdcox et al*, 1990]. *Nash and Sutcliffe* [1970] defined the coefficient of efficiency which ranges from minus infinity to 1.0, with higher values indicating better agreement, as

Mathematical symbols
not recognised

Columns not
recognised

[...] and is given by (Equation 1)
where x denotes...

Phase 3: Analysis of the corpus



Phase 3: Quantitative analysis of SCCERA

- ❖ Corpus analysis using WordSmith Tools 6.0 (Scott 2011)
- ❖ Comparisons across (a) RAs, (b) sub-sections, (c) sub-disciplines
- ❖ Word frequencies, keywords, key keywords, 2 to 8-word lexical bundles, type/token ratios, pedagogically significant concordance lines - e.g. disambiguation of near-synonymous words (Lee & Swales 2006)

Phase 3: Qualitative analysis of SCCERA

- ❖ Discourse analytical approach, investigating rhetorical characteristics of civil engineering RAs
- ❖ Move sequences in RA abstracts, introduction & discussion sections (often the most complex & problematic sections)
- ❖ Multimodality in civil engineering RAs

Word frequency (position)

et. al (24/5)

model (29)

fig (32)

we (35)

between (36)

time (37)

used (39)

results (44)

equation (45)

using (46)

table (52)

figure (59)

may (61)

values (64)

level (69)

analysis (72)

surface (76)

number (77)

study (82)

value (83)

models (89)

flow (93)

shown (92)

if (94)

case (95)

large (97)

project (98)

area (100)

effect (102)

due (104)

concrete (108)

method (109)

effects (112)

mean (113)

average (114)

same (115)

stress (116)

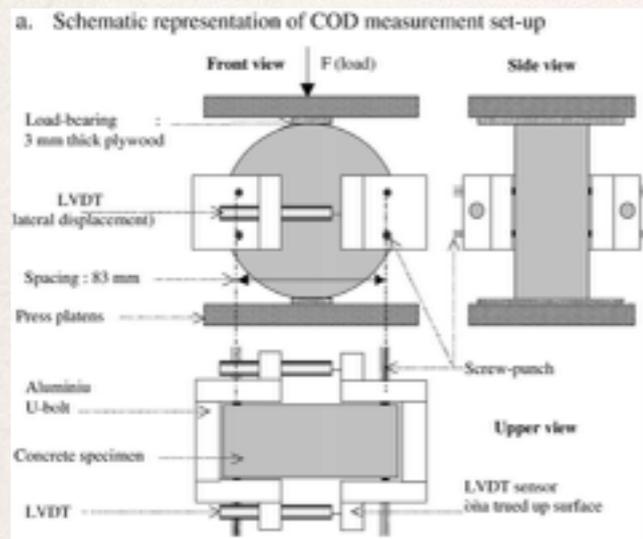
observed (117)

change (126)

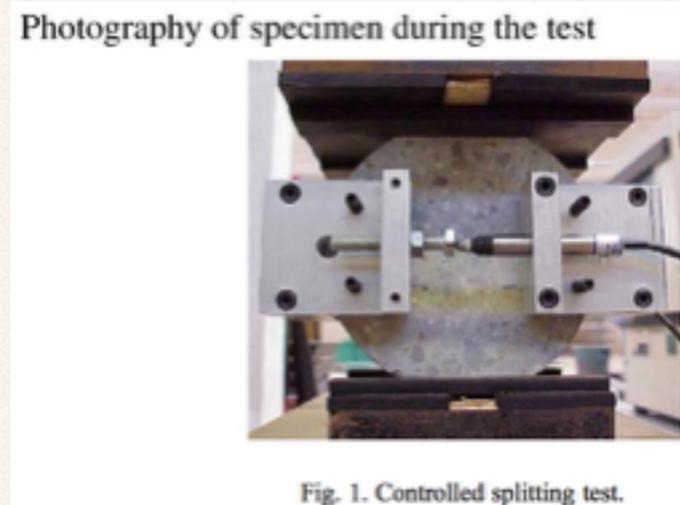
see (192)

Multimodality in civil engineering RAs

(see Fig. 1a)



(see Fig. 1)



[...], as in Eqn. (1):

$$J(x) = -D_e \frac{\partial c}{\partial x} + D_e \frac{zFE}{RTL} c + cv(x) \quad (1)$$

Fig. 3 presents...

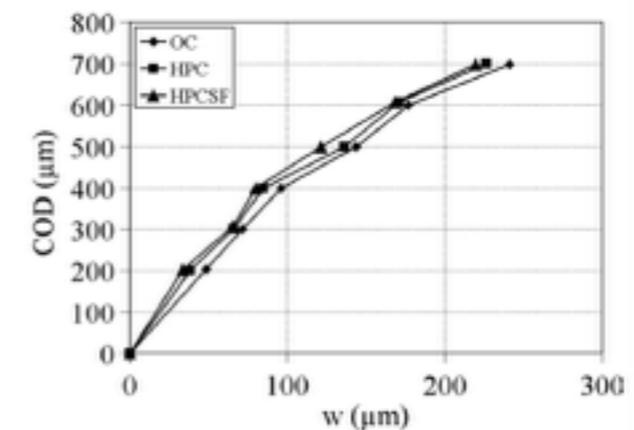
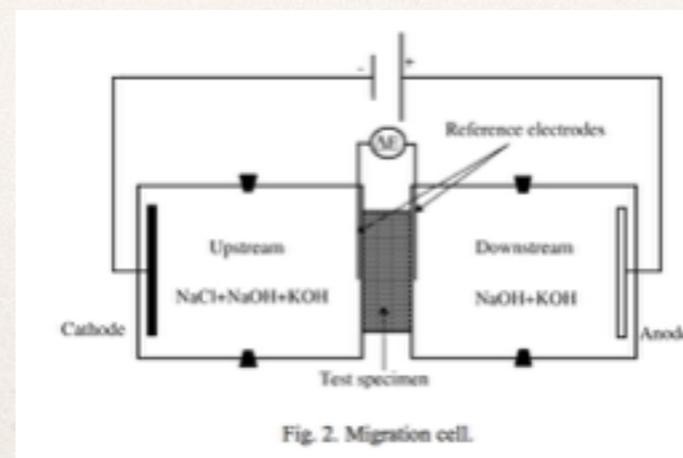


Fig. 3. Crack width versus crack opening displacement under loading.

(see Table 1)

Mix ingredients (kg/m ³)	OC	HPC	HPCSF
Coarse aggregate, 12.5–20 mm	777	550	579
Medium aggregate, 4–12.5 mm	415	475	465
Sand (Boulonnais), 0–5 mm	372	407	442
Sand (Seine), 0–4 mm	372	401	435
Cement CPA-CEM I 52.5	353	461	360
Silica fume	–	0	22
SP (e.s.)	–	12.4	12
Retarder (e.s.)	–	3.3	2.5
Total water	172	146	136
w/c	0.49	0.32	0.38
w/(C+SF)	0.49	0.32	0.36

(Fig. 2)



Epistemic language

Modal Verbs Could Couldn't May Might Should Shouldn't Would

Wouldn't Will Won't **Adjectives** Always Apparent Certain A certain

extent Clear Evident Possible Probable **Nouns** Claim Doubt Estimate

Evidence Possibility **Lexical Verbs** Appear Argue Assume

Believe Claim Doubt Estimate Expect Indicate Know Predict Presume Propose

Seem Speculate Suggest Suppose Tend Think **Adverbs** About Actually

Almost Apparently Approximately Around Certainly Clearly Definitely

Doubtless Essentially Evidently Frequently Generally In fact Indeed

Largely Likely Never Normally Obviously Of course Often Perhaps Possibly

Presumably Probably Quite Rarely Relatively Sometimes Surely Undoubtedly

Usually

Epistemic language

Modal Verbs Could Couldn't May Might Should Shouldn't Would

Wouldn't Will Won't **Adjectives** Always Apparent Certain A certain

extent **“Modes of knowing”:** **Nouns** Claim Doubt Estimate

Evidence Possibility **Lexical Verbs** Appear Argue Assume

Believe Claim Doubt Estimate Expect Indicate Know Predict Presume Propose

Seem Speculate Suggest Suppose Tend think **Adverbs** About Actually

Almost Apparently Approximately Around Certainly Clearly Definitely

Doubtless Essentially Evidently Frequently Generally In fact Indeed

Largely Likely Never Normally Obviously Of course Often Perhaps Possibly

Presumably Probably Quite Rarely Relatively Sometimes Surely Undoubtedly

Usually

Epistemic language

Modal Verbs Could Couldn't May Might Should Shouldn't Would

Wouldn't Will Won't **Adjectives** Always Apparent Certain A certain

extent Clear Evident Possible **Nouns** Claim Doubt Estimate
Evidence Possibility **Lexical Verbs** Appear Argue Assume

“Modes of knowing”:
Believe Claim Doubt Estimate Expect Indicate Know Predict Presume Propose
Seem Speculate Suggest Suppose Tend think **Adverbs** About Actually

Almost Apparently Approximately Around Certainly Clearly Definitely

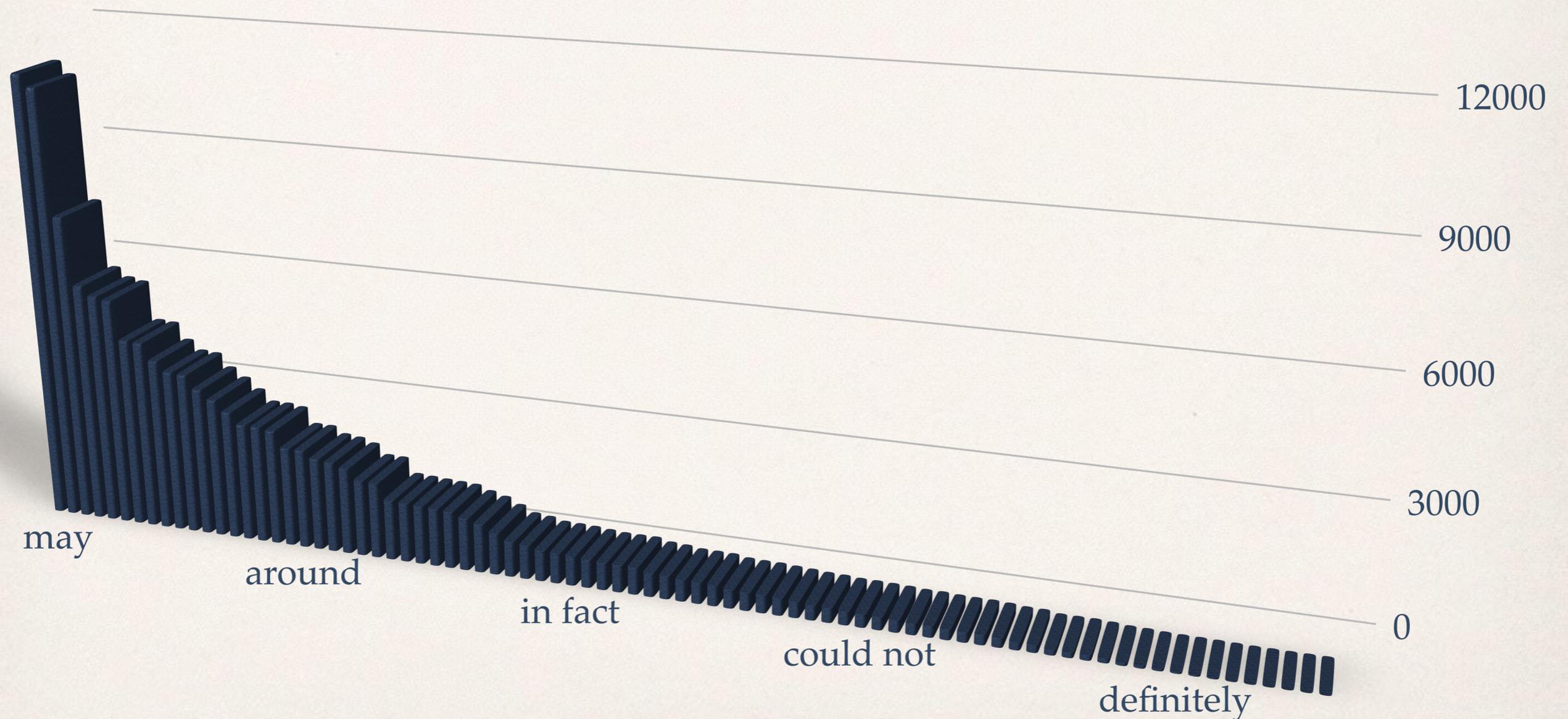
Doubtless Essentially Evidently Frequently Generally In fact Indeed
Largely Likely Never Normally Obviously Of course Often Perhaps Possibly

Presumably Probably Quite Rarely Relatively Sometimes Surely Undoubtedly

Usually

Most frequent epistemic items in academic writing (Hyland & Milton 1997)

Most frequent epistemic items in academic writing (Hyland & Milton 1997)



Epistemic items in SCCERA (frequency)

may (11,127)	could (4,044)	appear (2,267)	certain (1,204)
estimate(s) (10,823)	possible (3,713)	approximately (2,160)	quite (1,041)
will (7,703)	expect (3,488)	evidence (1,902)	argue (875)
about (6,019)	predict (3,223)	might (1,889)	indeed (833)
indicate(s) (5,794)	estimate (N) (2,942)	tend (1,502)	apparent (757)
would (5,722)	likely (2,934)	clear (1,478)	wouldn't (9)
should (4,754)	relatively (2,884)	seem (1,456)	won't (6)
assume (4,727)	often (2,486)	usually (1,439)	couldn't (3)
suggest (4,315)	around (2,480)	almost (1,419)	doubtless (3)
propose (4,074)	generally (2,311)	clearly (1,302)	shouldn't (2)

Modal expressions

will not

would

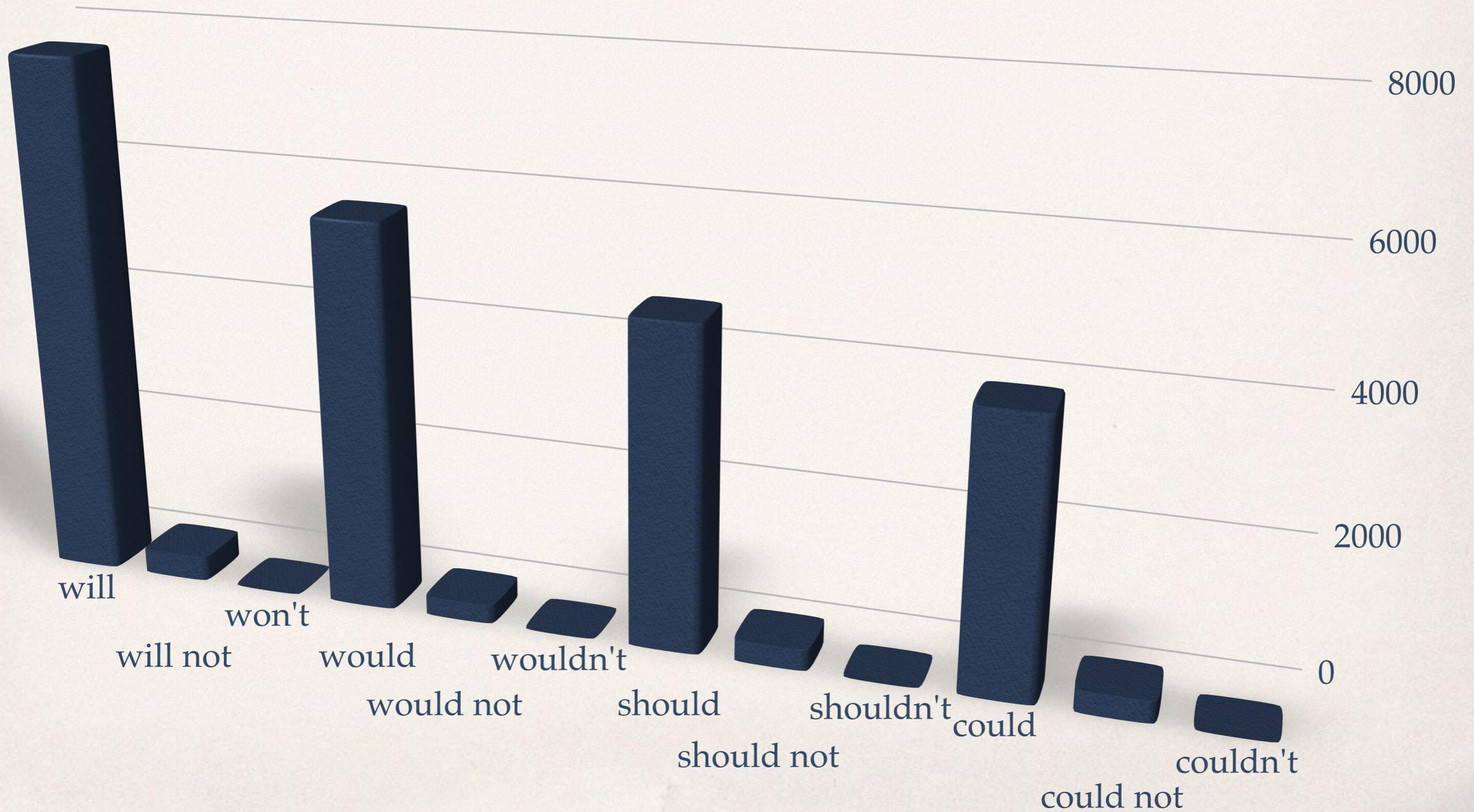
would not

should

should not

could not

Modal expressions



POS - CLAWS tagset

Position	POS tag	Info.	SCCERA	Medical	Brown
1	N	nouns	32.2%	29.1%	23.1%
2	V	verbs	13.4%	11.1%	15.5%
3	I	prepositions	13.4%	—	—
4	J	adjectives	10.2%	9.7%	6.9%

Coastal Engineering: Keywords vs. SCCERA (position)

wave (1)	storm (13)	numerical (23)	reef (33)
sea (2)	shoreline (14)	height (24)	waters (34)
coastal (3)	coast (15)	bed (25)	salinity (35)
ice (4)	erosion (16)	islands (26)	breakwater (37)
waves (5)	tidal (17)	water (27)	surge (38)
ocean (6)	tide (18)	shore (28)	swash (39)
breaking (7)	beaches (19)	offshore (29)	Atlantic (40)
beach (9)	currents (20)	island (30)	coasts (41)
shelf (10)	depth (21)	dune (31)	figure (42)
wind (12)	arctic (22)	runup (32)	shelves (44)

Coastal Engineering: Keywords vs. BNC (position)

wave (2)	ocean (14)	breaking (25)	tsunami (36)
et al (3/4)	water (16)	wind (26)	eq (37)
coastal (5)	figure (17)	tidal (27)	boundary (38)
ice (6)	velocity (18)	flow (28)	erosion (39)
model (8)	surface (19)	beach (29)	the (40)
fig (9)	depth (20)	height (30)	values (41)
equation (10)	numerical (21)	storm (32)	elevation (42)
sea (11)	sediment (22)	level (33)	measurements (43)
waves (12)	shelf (23)	measured (34)	salinity (45)
data (13)	shoreline (24)	results (35)	simulation (46)

Keywords: Hard vs. soft sub-disciplines of civil engineering

Hard (Mechanics & Structures)		Soft (Infra-structure dvlpt)	
damping	load	project	risk
response	steel	construction	pavement
beam	force	management	risks
structural	displacement	projects	team
stiffness	strain	cost	research
bridge	equation	success	safety
control	damage	life	leadership
vibration	frequency	costs	performance
damper	excitation	managers	process